

Group sparse optimization via $\ell_{p,q}$ regularization

Yaohua Hu^{*} Chong Li[†] Kaiwen Meng[‡] Jing Qin[§] Xiaoqi Yang[¶]

Abstract In this paper, we investigate a group sparse optimization problem via $\ell_{p,q}$ regularization in three aspects: theory, algorithm and application. In the theoretical aspect, by introducing a notion of group restricted eigenvalue condition, we establish some oracle property and a global recovery bound of order $O(\lambda^{\frac{2}{2-q}})$ for any point in a level set of the $\ell_{p,q}$ regularization problem, and by virtue of modern variational analysis techniques, we also provide a local analysis of recovery bound of order $O(\lambda^2)$ for a path of local minima. In the algorithmic aspect, we apply the well-known proximal gradient method to solve the $\ell_{p,q}$ regularization problems, either by analytically solving some specific $\ell_{p,q}$ regularization subproblems, or by using the Newton method to solve general $\ell_{p,q}$ regularization subproblems. In particular, we establish the linear convergence rate of the proximal gradient method for solving the $\ell_{1,q}$ regularization problem under some mild conditions. As a consequence, the linear convergence rate of proximal gradient method for solving the usual ℓ_q regularization problem ($0 < q < 1$) is obtained. Finally in the aspect of application, we present some numerical results on both the simulated data and the real data in gene transcriptional regulation.

Key words Group sparse optimization, $\ell_{p,q}$ regularization, nonconvex optimization, restricted eigenvalue condition, proximal gradient method, iterative thresholding algorithm, gene regulation network.

1 Introduction

In recent years, a great amount of attention has been paid to the sparse optimization problem, which is to find the sparse solutions of an underdetermined linear system. The sparse optimization problem arises in a wide range of fields, such as variable selection, pattern analysis, graphical modeling and compressive sensing; see [6, 11, 14, 20, 23, 48] and references therein.

^{*}College of Mathematics and Statistics, Shenzhen University, Shenzhen 518060, P. R. China (hyh19840428@163.com).

[†]Department of Mathematics, Zhejiang University, Hangzhou 310027, P. R. China (cli@zju.edu.cn).

[‡]School of Economics and Management, Southwest Jiaotong University, Chengdu 610031, P. R. China (mkwfly@126.com).

[§]School of Life Sciences, The Chinese University of Hong Kong, Shatin, New Territories, Hong Kong (qinjing@cuhk.edu.hk).

[¶]Department of Applied Mathematics, The Hong Kong Polytechnic University, Kowloon, Hong Kong (mayangxq@polyu.edu.hk).

In many applications, the underlying data usually can be represented approximately by a linear system of the form

$$Ax = b + \varepsilon,$$

where $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$ are known, $\varepsilon \in \mathbb{R}^m$ is an unknown noise vector, and $x = (x_1, x_2, \dots, x_n)^\top \in \mathbb{R}^n$ is the variable to be estimated. If $m \ll n$, the above linear system is seriously ill-conditioned and may have infinitely many solutions. The sparse optimization problem is to recover x from information b such that x is of a sparse structure. The sparsity of variable x has been measured by the ℓ_p norm $\|x\|_p$ ($p = 0$, see [6, 8]; $p = 1$, see [3, 14, 18, 20, 48, 53, 59]; and $p = 1/2$, see [13, 57]). The ℓ_p norm $\|x\|_p$ for $p > 0$ is defined as

$$\|x\|_p := \left(\sum_{i=1}^n |x_i|^p \right)^{1/p},$$

while the ℓ_0 norm $\|x\|_0$ is defined as the number of nonzero components of x . The sparse optimization problem can be modeled as

$$\begin{aligned} \min \quad & \|Ax - b\|_2 \\ \text{s.t.} \quad & \|x\|_0 \leq s, \end{aligned}$$

where s is the given sparsity level.

For the sparse optimization problem, a popular and practical technique is the regularization method, which is to transform the sparse optimization problem into an unconstrained optimization problem, called the regularization problem. For example, the ℓ_0 regularization problem is

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2 + \lambda \|x\|_0,$$

where $\lambda > 0$ is the regularization parameter, providing a tradeoff between accuracy and sparsity. However, the ℓ_0 regularization problem is nonconvex and non-Lipschitz, and thus it is generally intractable to solve it directly (indeed, it is NP-hard; see [38]).

To overcome this difficulty, two typical relaxations of the ℓ_0 regularization problem are introduced, which are the ℓ_1 regularization problem

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2 + \lambda \|x\|_1 \tag{1.1}$$

and the $\ell_{1/2}$ regularization problem

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2 + \lambda \|x\|_{1/2}^{1/2}. \tag{1.2}$$

1.1 ℓ_p regularization problems

The ℓ_1 regularization problem, also called Lasso [48] or Basis Pursuit [14], has attracted much attention and has been accepted as one of the most useful tools for the sparse optimization problem. Since the ℓ_1 regularization problem is a convex optimization problem,

many exclusive and efficient algorithms have been proposed and developed for solving (1.1), for instance, the interior-point methods [11, 14], LARs [22], the gradient projection method [25] and the alternating direction method [59]. However, in many practical applications, the solutions obtained from the ℓ_1 regularization problem are much less sparse than those of the ℓ_0 regularization problem, and it often leads to sub-optimal sparsity in reality; see, e.g., [13, 57, 63].

Recently, the $\ell_{1/2}$ regularization problem is proposed to improve the performance of sparsity recovery of the ℓ_1 regularization problem. Extensive computational studies in [13, 57] revealed that the $\ell_{1/2}$ regularization problem admits a significantly stronger sparsity promoting property than the ℓ_1 regularization problem in the sense that it guarantees to achieve the sparse solution from a smaller amount of samples. However, the $\ell_{1/2}$ regularization problem is nonconvex, nonsmooth and non-Lipschitz, and thus it is very difficult in general to design efficient algorithms for solving it. It was presented in [27] that finding the global minimal value of the $\ell_{1/2}$ regularization problem (1.2) is strongly NP-hard, while fortunately, computing a local minimum could be done in polynomial time. Some fast and efficient algorithms have been proposed to find a local minimum of (1.2), such as the hybrid OMP-SG algorithm [15] and the interior-point potential reduction algorithm [27].

The $\ell_{1/2}$ regularization problem (1.2) is a variant of lower-order penalty problems, investigated in [32, 33, 60], for a constrained optimization problem. The main advantage of the lower-order penalty functions over the classical ℓ_1 penalty functions is that they require weaker conditions to guarantee an exact penalization property and that their least exact penalty parameter is smaller; see [32]. It was reported in [60] that the first- and second-order necessary optimality conditions of lower order penalty problems converge to that of the original constrained optimization problem under a linearly independent constraint qualification.

Besides the preceding numerical algorithms, one of the most widely studied methods for solving the sparse optimization problem is the class of the iterative thresholding algorithms, which is studied in a uniform framework of proximal gradient methods; see [3, 6, 8, 16, 18, 39, 57] and references therein. It is convergent and of very low computational complexity. Benefitting from its simple formulation and low storage requirement, it is very efficient and applicable for large-scale sparse optimization problem. In particular, the iterative hard (resp. soft, half) thresholding algorithm for the ℓ_0 (resp. ℓ_1 , $\ell_{1/2}$) regularization problem was studied in [6, 8] (resp. [3, 18], [57]).

1.2 Global recovery bound

To estimate how far is the solution of regularization problems from that of the linear system, the global recovery bound or ℓ_2 consistency of the ℓ_1 regularization problem have been investigated in the literature [5, 9, 36, 51, 62]. More specifically, under some mild conditions on A , such as the restricted isometry property (RIP) [12] or restricted eigenvalue condition (REC) [5], van de Geer and Bühlmann [51] established a deterministic recovery bound for

the (convex) ℓ_1 regularization problem:

$$\|x^*(\ell_1) - \bar{x}\|_2^2 = O(\lambda^2 s), \quad (1.3)$$

where $x^*(\ell_1)$ is the solution of (1.1), \bar{x} is the solution of linear system $Ax = b$, and $s := \|\bar{x}\|_0$ is the sparsity of \bar{x} . In the statistic literature, [5, 9, 36, 62] provided the recovery bound in a high probability for the ℓ_1 regularization problem when the size of the variable tends to infinity, under REC/RIP or some relevant conditions. However, to the best of our knowledge, the recovery bound for the general (nonconvex) ℓ_p regularization problem is still undiscovered. We will establish such a deterministic property in Section 2.

1.3 Group sparse optimization

In applications, a wide class of problems usually has certain special structures, and recently, enhancing the recoverability due to the special structures has become an active topic in the sparse optimization. One of the most popular structures is the group sparsity structure, that is, the solution has a natural grouping of its components, and the components within each group are likely to be either all zeros or all nonzeros. In general, the grouping information can be any arbitrary partition of x , and it is usually pre-defined based on prior knowledge of specific problems. Let $x := (x_{\mathcal{G}_1}^\top, \dots, x_{\mathcal{G}_r}^\top)^\top$ represent the group structure of x . The group sparsity of x with such a group structure can be measured by an $\ell_{p,q}$ norm defined by

$$\|x\|_{p,q} := \left(\sum_{i=1}^r \|x_{\mathcal{G}_i}\|_p^q \right)^{1/q}.$$

Exploiting the group sparsity structure can reduce the degrees of freedom in the solution, thereby leading to better recovery performance. Benefitting from these advantages, the group sparse optimization model has been applied in birthweight prediction [2, 61], dynamic MRI [50] and gene finding [35, 58] with the $\ell_{2,1}$ norm. More specifically, the following $\ell_{2,1}$ regularization problem

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2^2 + \lambda \|x\|_{2,1}$$

was introduced by Yuan and Lin [61] to study the grouped variable selection in statistics under the name of group Lasso. The $\ell_{2,1}$ regularization, an important extension of the ℓ_1 regularization, proposes an ℓ_2 regularization for each group and ultimately yields the sparsity in the group manner. Since the $\ell_{2,1}$ regularization problem is a convex optimization problem, some effective algorithms have been proposed, such as, the spectral projected gradient method [52], SpaRSA [53] and the alternating direction method [19].

1.4 The aim of the paper

In this paper, we will investigate the group sparse optimization via $\ell_{p,q}$ ($p \geq 1, 0 \leq q \leq 1$) regularization, also called the $\ell_{p,q}$ regularization problem

$$\min_{x \in \mathbb{R}^n} F(x) := \|Ax - b\|_2^2 + \lambda \|x\|_{p,q}^q. \quad (1.4)$$

We will investigate the oracle property and recovery bound for the $\ell_{p,q}$ regularization problem, which extend the existing results in two ways: one is the lower-order regularization problem, including the ℓ_q ($q < 1$) regularization problem; the other is the group sparse optimization problem, including the $\ell_{2,1}$ regularization problem (group Lasso) as a special case. To this end, we introduce the weaker notions of REC: the lower-order REC and the group REC (GREC). We will further establish the relationships between the new notions with the classical one: the lower-order REC is weaker than the classical REC, but the reverse is not true (see Example 2.1); and the GREC is weaker than the REC. Under the GREC, we will provide the oracle property and the global recovery bound for the $\ell_{p,q}$ regularization problem (see Theorem 2.1). Furthermore, we will conduct a local analysis of recovery bound for the $\ell_{p,q}$ regularization problem by virtue of modern variational analysis techniques [45]. More precisely, we assume that the columns of A corresponding to the active components of \bar{x} (a solution of $Ax = b$) are linearly independent. This leads us to the application of implicit function theorem and thus guarantees the existence of a local path around \bar{x} which satisfies a second-order growth condition. As such, in the local recovery bound, we establish a uniform quadratic recovery bound for all $\ell_{p,q}$ regularization problem, see Theorem 2.2.

The proximal gradient method is one of the most popular and practical methods for the sparse optimization problems, either convex or nonconvex problems. We will apply the proximal gradient method to solve the $\ell_{p,q}$ regularization problem (1.4). The advantage of the use of the proximal gradient method is that for some specific regularization problems, the proximal subproblems have the analytical solutions, and the resulting algorithm is thus practically attractive. In the general cases when the analytical solutions of the proximal optimization subproblems seem not available, we will employ the Newton method to solve the proximal optimization subproblems. Furthermore, we will investigate the linear convergence rate of proximal gradient method for solving the $\ell_{p,q}$ regularization problem when $p = 1$ and $0 < q < 1$ under the assumption that any nonzero group of a local minimum is active. The problem (1.4) of the case $p = 1$ and $0 < q < 1$ possesses the properties that the regularization term $\|x\|_{p,q}^q$ is concave and the objective function $F(x)$ of (1.4) satisfies a second-order growth condition, which play an important role in the establishment of the linear convergence rate. To the best of our knowledge, this is the first attempt to study the linear convergence rate of proximal gradient method for solving nonconvex optimization problems. As a consequence of this result, we will obtain the linear convergence rate of proximal gradient method for solving ℓ_q regularization problem ($0 < q < 1$), which includes the iterative half thresholding algorithm ($q = 1/2$) proposed in [57] as a special case. The result on linear convergence rate of proximal gradient method for solving ℓ_q regularization problem is still new, as far as we know.

In the aspect of application, we will conduct some numerical experiments on both simulated data and real data in gene transcriptional regulation to demonstrate the performance of the proposed proximal gradient method. From the numerical results, it is demonstrated that the $\ell_{p,1/2}$ regularization is the best one among the $\ell_{p,q}$ regularizations for $q \in [0, 1]$, and it outperforms the $\ell_{p,1}$ and $\ell_{p,0}$ regularizations on both accuracy and robustness. This obser-

vation is consistent with several previous numerical studies [13, 57] on the ℓ_p regularization problem.

1.5 Main contributions

The main objectives of this paper are to establish the oracle property and recovery bound, to design an efficient numerical method for the $\ell_{p,q}$ regularization problem (1.4), and to apply the proposed method to the gene transcriptional regulation. The main contributions are presented as follows.

- (i) We establish the following global recovery bound for the $\ell_{p,q}$ regularization problem (1.4) under the (p, q) -GREC:

$$\|x^* - \bar{x}\|_2^2 \leq \begin{cases} O\left(\lambda^{\frac{2}{2-q}} S\right), & 2^{K-1}q = 1, \\ O\left(\lambda^{\frac{2}{2-q}} S^{\frac{3-q}{2-q}}\right), & 2^{K-1}q > 1, \end{cases} \quad (1.5)$$

where \bar{x} is a true solution of $Ax = b$, $S := \|\bar{x}\|_{p,0}$ is the group sparsity of \bar{x} , $0 < q \leq 1 \leq p \leq 2$, x^* is any point in the level set $\text{lev}_F(\bar{x})$ of (1.4), and K is the smallest integer such that $2^{K-1}q \geq 1$.

- (ii) By virtue of the variational analysis technique, for all the $\ell_{p,q}$ regularization problems, we establish a uniform local recovery bound

$$\|x_{p,q}^*(\lambda) - \bar{x}\|_2^2 \leq O(\lambda^2 S) \quad \text{for small } \lambda,$$

where $0 < q < 1 \leq p$ and $x_{p,q}^*(\lambda)$ is a local optimal solution of (1.4) (near \bar{x}).

- (iii) We present the analytical formulae for the proximal optimization subproblems of some specific $\ell_{p,q}$ regularizations, when $p = 1, 2$ and $q = 0, 1/2, 2/3, 1$. Moreover, we prove that any sequence $\{x^k\}$, generated by proximal gradient method for solving the $\ell_{1,q}$ regularization problem, linearly converges to a local minimum x^* under some mild conditions, i.e., there exist $N \in \mathbb{N}$, $C > 0$ and $\eta \in (0, 1)$ such that

$$F(x^k) - F(x^*) \leq C\eta^k \quad \text{and} \quad \|x^k - x^*\|_2 \leq C\eta^k, \quad \text{for any } k \geq N.$$

- (iv) Our numerical experiments show that, measured by the biological golden standards, the accuracy of the gene regulation networks forecasting can be improved by exploiting the group structure of TF complexes. The successful application of group sparse optimization to gene transcriptional regulation will facilitate biologists to study the gene regulation of higher model organisms in a genome-wide scale.

1.6 The organization of the paper

This paper is organized as follows. In section 2, we introduce the concepts of q -REC and GREC, and establish the oracle property and (global and local) recovery bounds for the $\ell_{p,q}$ regularization problem. In section 3, we apply the proximal gradient method to solve the group sparse optimization using different types of $\ell_{p,q}$ regularization, and investigate the linear convergence rate of the resulting proximal gradient method. Finally, section 4 exhibits the numerical results on both simulated data and real data in gene transcriptional regulation.

2 Global and local recovery bounds

This section is devoted to the study of the oracle property and (global and local) recovery bounds for the $\ell_{p,q}$ regularization problem (1.4). To this end, we first present some basic inequalities of ℓ_p norm and introduce the notions of RECs, as well as their relationships.

The notation adopted in this paper is described as follows. We let the lowercase letters x, y, z denote the vectors, capital letters N, S denote the numbers of groups in the index sets, caligraphic letters $\mathcal{I}, \mathcal{T}, \mathcal{S}, \mathcal{J}, \mathcal{N}$ denote the index sets. In particular, we use \mathcal{G}_i to denote the index set corresponding to the i -th group and $\mathcal{G}_{\mathcal{S}}$ to denote the index set $\{\mathcal{G}_i : i \in \mathcal{S}\}$. For $x \in \mathbb{R}^n$ and $\mathcal{T} \subseteq \{1, \dots, n\}$, we use $x_{\mathcal{T}}$ to denote the subvector of x corresponding to \mathcal{T} .

Throughout this paper, we assume that the group sparse optimization problem is of the group structure described as follows. Let $x := (x_{\mathcal{G}_1}^\top, \dots, x_{\mathcal{G}_r}^\top)^\top$ represent the group structure of x , where $\{x_{\mathcal{G}_i} \in \mathbb{R}^{n_i} : i = 1, \dots, r\}$ is the grouping of x , $\sum_{i=1}^r n_i = n$ and $n_{\max} := \max\{n_i : i \in \{1, \dots, r\}\}$. For a group $x_{\mathcal{G}_i}$, we use $x_{\mathcal{G}_i} = 0$ (reps. $x_{\mathcal{G}_i} \neq 0$, $x_{\mathcal{G}_i} \neq_{\mathbf{a}} 0$) to denote a zero (reps. nonzero, active) group, where $x_{\mathcal{G}_i} = 0$ means that $x_j = 0$ for all $j \in \mathcal{G}_i$; $x_{\mathcal{G}_i} \neq 0$ means that $x_j \neq 0$ for some $j \in \mathcal{G}_i$; and $x_{\mathcal{G}_i} \neq_{\mathbf{a}} 0$ means that $x_j \neq 0$ for all $j \in \mathcal{G}_i$. It is trivial to see that

$$x_{\mathcal{G}_i} \neq_{\mathbf{a}} 0 \quad \Rightarrow \quad x_{\mathcal{G}_i} \neq 0.$$

For this group structure and $p > 0$, the $\ell_{p,q}$ norm of x is defined by

$$\|x\|_{p,q} = \begin{cases} (\sum_{i=1}^r \|x_{\mathcal{G}_i}\|_p^q)^{1/q}, & q > 0, \\ \sum_{i=1}^r \|x_{\mathcal{G}_i}\|_p^0, & q = 0, \end{cases} \quad (2.1)$$

which proposes the ℓ_p norm for each group and then processes the ℓ_q norm for the resulting vector. When $p = q$, the $\ell_{p,q}$ norm coincides with the ℓ_p norm, i.e., $\|x\|_{p,p} = \|x\|_p$. Furthermore, all $\ell_{p,0}$ norms share the same formula, i.e., $\|x\|_{p,0} = \|x\|_{2,0}$, for all $p > 0$. In particular, when the grouping structure is degenerated to the individual feature level, i.e., $n_{\max} = 1$ or $n = r$, we have $\|x\|_{p,q} = \|x\|_q$ for all $p > 0$ and $q > 0$.

Moreover, we assume that A and b in (1.4) are related by a linear model (noiseless)

$$b = A\bar{x}.$$

Let $\mathcal{S} := \{i \in \{1, \dots, r\} : \bar{x}_{\mathcal{G}_i} \neq 0\}$ be the index set of nonzero groups of \bar{x} , $\mathcal{S}^c := \{1, \dots, r\} \setminus \mathcal{S}$ be the complement of \mathcal{S} , $S := |\mathcal{S}|$ be the group sparsity of \bar{x} , and $n_{\mathbf{a}} := \sum_{i \in \mathcal{S}} n_i$.

2.1 Inequalities of $\ell_{p,q}$ norm

We begin with some basic inequalities of ℓ_p and $\ell_{p,q}$ norms, which will be useful in the later discussion of RECs and recovery bounds. First, we recall the following well-known inequality

$$\left(\sum_{i=1}^n |x_i|^{\gamma_2} \right)^{1/\gamma_2} \leq \left(\sum_{i=1}^n |x_i|^{\gamma_1} \right)^{1/\gamma_1} \quad \text{if } 0 < \gamma_1 \leq \gamma_2, \quad (2.2)$$

or equivalently $(x = (x_1, x_2, \dots, x_n)^\top)$,

$$\|x\|_{\gamma_2} \leq \|x\|_{\gamma_1} \quad \text{if } 0 < \gamma_1 \leq \gamma_2.$$

The following lemma improves [32, Lemma 4.1] and extends to the $\ell_{p,q}$ norm. It will be useful in providing a shaper global recovery bound (see Theorem 2.1 below).

Lemma 2.1. *Let $0 < q \leq p \leq 2$, $x \in \mathbb{R}^n$ and K be the smallest integer such that $2^{K-1}q \geq 1$. Then the following relations hold.*

- (i) $\|x\|_q^q \leq n^{1-2^{-K}} \|x\|_2^q.$
- (ii) $\|x\|_{p,q}^q \leq n^{1-2^{-K}} \|x\|_{p,2}^q.$

Proof. (i) Repeatedly using the property that $\|x\|_1 \leq \sqrt{n}\|x\|_2$, one has that

$$\begin{aligned} \|x\|_q^q &\leq \sqrt{n} \left(\sum_{i=1}^n |x_i|^{2q} \right)^{1/2} \\ &\leq \dots \\ &\leq n^{\frac{1}{2} + \dots + \frac{1}{2^K}} \left(\sum_{i=1}^n |x_i|^{2^K q} \right)^{2^{-K}} \end{aligned}$$

Since $2^{K-1}q \geq 1$, by (2.2), we have

$$\left(\sum_{i=1}^n |x_i|^{2^K q} \right)^{2^{-K}} = \left(\sum_{i=1}^n (|x_i|^2)^{2^{K-1}q} \right)^{\frac{1}{2^{K-1}q} \frac{q}{2}} \leq \left(\sum_{i=1}^n |x_i|^2 \right)^{q/2} = \|x\|_2^q.$$

Therefore, we arrive at the conclusion that

$$\|x\|_q^q \leq n^{1-2^{-K}} \|x\|_2^q.$$

- (ii) By (2.1), it is a direct consequence of (i). ■

For example, if $q = 1$, then $K = 1$; if $q = \frac{1}{2}$ or $\frac{2}{3}$, then $K = 2$. The following lemma describes the triangle inequality of $\|\cdot\|_{p,q}^q$.

Lemma 2.2. *Let $0 < q \leq 1 \leq p$ and $x, y \in \mathbb{R}^n$. Then*

$$\|x\|_{p,q}^q - \|y\|_{p,q}^q \leq \|x - y\|_{p,q}^q.$$

Proof. By the subadditivity of ℓ_p norm and (2.2), it is easy to see that

$$\|x_{\mathcal{G}_i}\|_p^q - \|y_{\mathcal{G}_i}\|_p^q \leq \|x_{\mathcal{G}_i} - y_{\mathcal{G}_i}\|_p^q, \quad \text{for } i = 1, \dots, r.$$

Consequently, the conclusion directly follows from (2.1). \blacksquare

The following lemma will be beneficial to studying properties of the lower-order REC in Proposition 2.1.

Lemma 2.3. *Let $\gamma \geq 1$, and two finite sequences $\{y_i : i \in \mathcal{I}\}$ and $\{x_j : j \in \mathcal{J}\}$ satisfy that $y_i \geq x_j \geq 0$ for all $(i, j) \in \mathcal{I} \times \mathcal{J}$. If $\sum_{i \in \mathcal{I}} y_i \geq \sum_{j \in \mathcal{J}} x_j$, then $\sum_{i \in \mathcal{I}} y_i^\gamma \geq \sum_{j \in \mathcal{J}} x_j^\gamma$.*

Proof. Set $\bar{y} := \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} y_i$ and $\alpha := \min_{i \in \mathcal{I}} y_i$. By [32, Lemma 4.1(ii)], one has that

$$\sum_{i \in \mathcal{I}} y_i^\gamma \geq \frac{1}{|\mathcal{I}|^{\gamma-1}} \left(\sum_{i \in \mathcal{I}} y_i \right)^\gamma = |\mathcal{I}| \bar{y}^\gamma. \quad (2.3)$$

On the other hand, let $M \in \mathbb{N}$ and $\beta \in [0, \alpha)$ be such that $\sum_{j \in \mathcal{J}} x_j = M\alpha + \beta$. Observing $\gamma \geq 1$ and $0 \leq x_j \leq \alpha$ for all $j \in \mathcal{J}$, we obtain that $x_j^\gamma \leq x_j \alpha^{\gamma-1}$, and thus, $\sum_{j \in \mathcal{J}} x_j^\gamma \leq M\alpha^\gamma + \alpha^{\gamma-1}\beta$. By (2.3), it remains to show that

$$|\mathcal{I}| \bar{y}^\gamma \geq M\alpha^\gamma + \alpha^{\gamma-1}\beta. \quad (2.4)$$

If $|\mathcal{I}| > M$, the relation (2.4) is trivial since $\bar{y} \geq \alpha > \beta$; otherwise, $|\mathcal{I}| \leq M$, from the facts that $|\mathcal{I}| \bar{y} \geq M\alpha + \beta$ (i.e., $\sum_{i \in \mathcal{I}} y_i \geq \sum_{j \in \mathcal{J}} x_j$) and that $\gamma \geq 1$, it follows that

$$|\mathcal{I}| \bar{y}^\gamma \geq M^{1-\gamma} (M\alpha + \beta)^\gamma \geq M^{1-\gamma} (M^\gamma \alpha^\gamma + \gamma M^{\gamma-1} \alpha^{\gamma-1} \beta) \geq M\alpha^\gamma + \alpha^{\gamma-1} \beta.$$

Therefore, we obtain the relation (2.4), and the proof is complete. \blacksquare

2.2 Group restricted eigenvalue conditions

This subsection aims at the development of the critical conditions on the matrix A to guarantee the oracle property and the global recovery bound of the $\ell_{p,q}$ regularization problem (1.4). In particular, we will focus on the restricted eigenvalue condition (REC), and extend it to the lower-order setting and equip it with the group structure.

In the scenario of sparse optimization, given the sparsity level s , it is always assumed that the $2s$ -sparse minimal eigenvalue of $A^\top A$ is positive (see, e.g., [5, 9, 36]), that is,

$$\phi_{\min}(2s) := \min_{\|x\|_0 \leq 2s} \frac{x^\top A^\top A x}{x^\top x} > 0, \quad (2.5)$$

which is the minimal eigenvalue of any $2s \times 2s$ dimensional submatrix. It is well-known that the solution at sparsity level s of the linear system $Ax = b$ is unique if the condition (2.5) is satisfied; otherwise, assume that there are two distinct vectors \hat{x} and \tilde{x} such that $A\hat{x} = A\tilde{x}$ and $\|\hat{x}\|_0 = \|\tilde{x}\|_0 = s$. Then such $x := \hat{x} - \tilde{x}$ is a vector so that $Ax = 0$ and $\|x\|_0 \leq 2s$,

and thus $\phi_{\min}(2s) = 0$, which is contradict with (2.5). Therefore, if the $2s$ -sparse minimal eigenvalue of $A^\top A$ is zero (i.e., $\phi_{\min}(2s) = 0$), one has no hope of recovering the true sparse solution from noisy observations.

However, only the condition (2.5) is not enough and some further condition is required to maintain the nice recovery of the regularization problem; see [5, 9, 36, 51, 62] and references therein. For example, the REC is introduced in Bickel et al. [5] to investigate the ℓ_2 consistency of the ℓ_1 regularization problem (Lasso), where the minimum in (2.5) is replaced by the minimum over a restricted set of vectors measured by an ℓ_1 norm inequality and the denominator is replaced by the ℓ_2 norm of only a part of x .

We now introduce the lower-order REC. Note that the residual $\hat{x} := x^*(\ell_q) - \bar{x}$, where $x^*(\ell_q)$ is an optimal solution of ℓ_q regularization problem and \bar{x} is a sparse solution of $Ax = b$, of the ℓ_q regularization problem always satisfies

$$\|\hat{x}_{\mathcal{S}^c}\|_q \leq \|\hat{x}_{\mathcal{S}}\|_q, \quad (2.6)$$

where \mathcal{S} is the support of \bar{x} . Thus we introduce a lower-order REC, where the minimum is taken over a restricted set measured by an ℓ_q norm inequality such as (2.6), for establishing the global recovery bound of the ℓ_q regularization problem. Given $s \leq t \ll n$, $x \in \mathbb{R}^n$ and $\mathcal{I} \subset \{1, \dots, n\}$, we denote by $\mathcal{I}(x; t)$ the subset of $\{1, \dots, n\}$ corresponding to the first t largest coordinates in absolute value of x in \mathcal{I}^c .

Definition 2.1. Let $0 \leq q \leq 1$. The q -restricted eigenvalue condition relative to (s, t) (q -REC(s, t)) is said to be satisfied if

$$\phi_q(s, t) := \min \left\{ \frac{\|Ax\|_2}{\|x_{\mathcal{T}}\|_2} : |\mathcal{I}| \leq s, \|x_{\mathcal{I}^c}\|_q \leq \|x_{\mathcal{I}}\|_q, \mathcal{T} = \mathcal{I}(x; t) \cup \mathcal{I} \right\} > 0.$$

The q -REC describes a kind of restricted positive definiteness of $A^\top A$, which is valid only for the vectors satisfying the relation measured by an ℓ_q norm. The q -REC presents a unified framework of the REC-type conditions when $q \in [0, 1]$. In particular, we note by definition that 1-REC reduces to the classical REC [5], and that $\phi_{\min}(2s) = \phi_0^2(s, s)$, and thus

$$(2.5) \Leftrightarrow 0\text{-REC}(s, s) \text{ is satisfied.}$$

It is well-known in the literature that the 1-REC is a stronger condition than the 0-REC (i.e., (2.5)). A natural question arises what are the relationship between the general q -RECs. To answer this question, associated with the q -REC, we consider the feasible set

$$C_q(s) := \{x \in \mathbb{R}^n : \|x_{\mathcal{I}^c}\|_q \leq \|x_{\mathcal{I}}\|_q \text{ for some } |\mathcal{I}| \leq s\},$$

which is a cone. Since the objective function associated with the q -REC is homogeneous, the q -REC(s, t) says that the null space of A does not cross over $C_q(s)$. Figure ?? presents the geometric interpretation of the q -RECs. It is shown that $C_0(s) \subseteq C_{1/2}(s) \subseteq C_1(s)$, and thus

$$1\text{-REC} \Rightarrow 1/2\text{-REC} \Rightarrow 0\text{-REC}.$$

It is also observed from Figure ?? that the gap between the 1-REC and 1/2-REC and that between 1/2-REC and 0-REC are the matrices whose null spaces all fall in the cones of $C_1(s) \setminus C_{1/2}(s)$ and $C_{1/2}(s) \setminus C_0(s)$, respectively.

We now provide a rigorous proof in the following proposition to identify the relationship between the feasible sets $C_q(s)$ and between the general q -RECs: the lower the q , the smaller the cone $C_q(s)$, and the weaker the q -REC.

Proposition 2.1. *Let $0 \leq q_1 \leq q_2 \leq 1$ and $1 \leq s \leq t \ll n$. Then the following statements are true:*

- (i) $C_{q_1}(s) \subseteq C_{q_2}(s)$, and
- (ii) if the q_2 -REC(s, t) holds, then the q_1 -REC(s, t) holds.

Proof. (i) Let $x \in C_{q_1}(s)$. We use \mathcal{I}_* to denote the index set of the first s largest coordinates in absolute value of x . Since $x \in C_{q_1}(s)$, it follows that $\|x_{\mathcal{I}_*^c}\|_{q_1} \leq \|x_{\mathcal{I}_*}\|_{q_1}$ ($|\mathcal{I}_*| \leq s$ due to the construction of \mathcal{I}_*). By Lemma 2.3 (taking $\gamma = q_2/q_1$), one has that

$$\|x_{\mathcal{I}_*^c}\|_{q_2} \leq \|x_{\mathcal{I}_*}\|_{q_2},$$

that is, $x \in C_{q_2}(s)$. Hence it follows that $C_{q_1}(s) \subseteq C_{q_2}(s)$.

- (ii) As proved by (i) that $C_{q_1}(s) \subseteq C_{q_2}(s)$, by the definition of q -REC, it follows that

$$\phi_{q_1}(s, t) \geq \phi_{q_2}(s, t) > 0.$$

The proof is complete. ■

To the best of our knowledge, this is the first work on introducing the lower-order REC and establishing the relationship of the lower-order RECs. In the following, we provide a counter example to show that the reverse of Proposition 2.1 is not true.

Example 2.1 (A matrix satisfying 1/2-REC but not REC). *Consider the matrix*

$$A = \begin{pmatrix} a & a+c & a-c \\ \tilde{a} & \tilde{a}-\tilde{c} & \tilde{a}+\tilde{c} \end{pmatrix} \in \mathbb{R}^{2 \times 3},$$

where $a > c > 0$ and $\tilde{a} > \tilde{c} > 0$. This matrix A does not satisfy the REC(1,1). Indeed, by letting $\mathcal{J} = \{1\}$ and $x = (2, -1, -1)^\top$, we have $Ax = 0$ and thus $\phi(1, 1) = 0$.

However, A satisfies the 1/2-REC(1,1). It suffices to show that $\phi_{1/2}(1, 1) > 0$. Let $x = (x_1, x_2, x_3)^\top$ satisfy the constraint associated with 1/2-REC(1,1). As $s = 1$, the deduction is divided into the following three cases.

(i) $\mathcal{J} = \{1\}$. Then

$$|x_1| \geq \|x_{\mathcal{J}^c}\|_{1/2} = |x_2| + |x_3| + 2|x_2|^{1/2}|x_3|^{1/2}. \quad (2.7)$$

Without loss of generality, we assume $|x_1| \geq |x_2| \geq |x_3|$. Hence, $\mathcal{T} = \{1, 2\}$ and

$$\frac{\|Ax\|_2}{\|x_{\mathcal{T}}\|_2} \geq \frac{\min\{a, \tilde{a}\}|x_1 + x_2 + x_3| + \min\{c, \tilde{c}\}|x_2 - x_3|}{|x_1| + |x_2|}. \quad (2.8)$$

If $|x_2| \leq \frac{1}{3}|x_1|$, (2.8) reduces to

$$\frac{\|Ax\|_2}{\|x_{\mathcal{T}}\|_2} \geq \frac{\frac{\min\{a, \tilde{a}\}}{3}|x_1|}{\frac{4}{3}|x_1|} = \frac{\min\{a, \tilde{a}\}}{4}. \quad (2.9)$$

If $|x_2| \geq \frac{1}{3}|x_1|$, substituting (2.7) into (2.8), one has that

$$\frac{\|Ax\|_2}{\|x_{\mathcal{T}}\|_2} \geq \frac{2 \min\{a, \tilde{a}\}|x_2|^{1/2}|x_3|^{1/2} + \min\{c, \tilde{c}\}|x_2 - x_3|}{4|x_2|} \geq \begin{cases} \frac{\min\{c, \tilde{c}\}}{8}, & |x_3| \leq \frac{1}{2}|x_2|, \\ \frac{\min\{a, \tilde{a}\}}{4}, & |x_3| \geq \frac{1}{2}|x_2|. \end{cases} \quad (2.10)$$

(ii) $\mathcal{J} = \{2\}$. Since $\mathcal{T} = \{2, 1\}$ or $\{2, 3\}$, it follows from [32, Lemma 4.1(ii)] that

$$|x_2| \geq \|x_{\mathcal{J}^c}\|_{1/2} \geq |x_1| + |x_3|. \quad (2.11)$$

Thus, it is easy to verify that $\|x_{\mathcal{T}}\|_2 \leq 2|x_2|$ and that

$$\frac{\|Ax\|_2}{\|x_{\mathcal{T}}\|_2} \geq \frac{|ax_1 + (a+c)x_2 + (a-c)x_3|}{2|x_2|} = \frac{|a(x_1 + x_2 + \frac{a-c}{a}x_3) + cx_2|}{2|x_2|} \geq \frac{c}{2}, \quad (2.12)$$

where the last inequality follows from (2.11) and the fact that $a > c$.

(iii) $\mathcal{J} = \{3\}$. Similar to the deduction of (ii), we have that

$$\frac{\|Ax\|_2}{\|x_{\mathcal{T}}\|_2} \geq \frac{|\tilde{a}x_1 + (\tilde{a} - \tilde{c})x_2 + (\tilde{a} + \tilde{c})x_3|}{2|x_3|} \geq \frac{\tilde{c}}{2}. \quad (2.13)$$

Therefore, by (2.9), (2.10), (2.12) and (2.13), one has that $\phi_{1/2}(1, 1) \geq \frac{1}{8} \min\{c, \tilde{c}\} > 0$, and thus, the matrix A satisfies the $1/2$ -REC(1,1).

In order to establish the oracle property and the global recovery bound for the $\ell_{p,q}$ regularization problem, we further introduce the notion of group restricted eigenvalue condition (GREC). Given $S \leq N \ll r$, $x \in \mathbb{R}^n$ and $\mathcal{J} \subset \{1, \dots, r\}$, we use $\text{rank}_i(x)$ to denote the rank of $\|x_{\mathcal{G}_i}\|_p$ among $\{\|x_{\mathcal{G}_j}\|_p : j \in \mathcal{J}^c\}$ (in a decreasing order), $\mathcal{J}(x; N)$ to denote the index set of the first N largest groups in the value of $\|x_{\mathcal{G}_i}\|_p$ among $\{\|x_{\mathcal{G}_j}\|_p : j \in \mathcal{J}^c\}$, that is,

$$\mathcal{J}(x; N) := \{i \in \mathcal{J}^c : \text{rank}_i(x) \in \{1, \dots, N\}\}.$$

Furthermore, by letting $R := \lceil \frac{r-|\mathcal{J}|}{N} \rceil$, we denote

$$\mathcal{J}_k(x; N) := \begin{cases} \{i \in \mathcal{J}^c : \text{rank}_i(x) \in \{kN + 1, \dots, (k+1)N\}\}, & k = 1, \dots, R-1, \\ \{i \in \mathcal{J}^c : \text{rank}_i(x) \in \{RN + 1, \dots, r - |\mathcal{J}|\}\}, & k = R. \end{cases} \quad (2.14)$$

Note that the residual $\hat{x} := x^*(\ell_{p,q}) - \bar{x}$ of the $\ell_{p,q}$ regularization problem always satisfies $\|\hat{x}_{\mathcal{G}_{S^c}}\|_{p,q} \leq \|\hat{x}_{\mathcal{G}_S}\|_{p,q}$. Thus we introduce the notion of GREC, where the minimum is taken over a restricted set measured by an $\ell_{p,q}$ norm inequality, as follows.

Definition 2.2. Let $0 < q \leq p \leq 2$. The (p, q) -group restricted eigenvalue condition relative to (S, N) ((p, q) -GREC(S, N)) is said to be satisfied if

$$\phi_{p,q}(S, N) := \min \left\{ \frac{\|Ax\|_2}{\|x_{\mathcal{G}_N}\|_{p,2}} : |\mathcal{J}| \leq S, \|x_{\mathcal{G}_{\mathcal{J}^c}}\|_{p,q} \leq \|x_{\mathcal{G}_S}\|_{p,q}, \mathcal{N} = \mathcal{J}(x; N) \cup \mathcal{J} \right\} > 0.$$

The (p, q) -GREC extends the q -REC to the setting equipping with a pre-defined group structure. Handling the components in each group as one element, the (p, q) -GREC admits the fewer degree of freedom, which is S , about s/n_{\max} , on its associated constraint than that of the q -REC, and thus it characterizes a weaker condition than the q -REC. For example, the 0-REC(s, s) is to indicate the restricted positive definiteness of $A^\top A$, which is valid only for the vectors whose cardinality is less than $2s$; while the $(p, 0)$ -GREC(S, S) is to describe the restricted positive definiteness of $A^\top A$ on any $2S$ -group support, whose degree of freedom is much less than the $2s$ -support. Thus the $(p, 0)$ -GREC(S, S) provides a broader condition than the 0-REC(s, s). Similar to the proof of Proposition 2.1, we can show that if $0 \leq q_1 \leq q_2 \leq 1 \leq p \leq 2$ and the (p, q_2) -GREC(S, N) holds, then the (p, q_1) -GREC(S, N) also holds.

We end this subsection by providing the following lemma, which will be useful in establishing the global recovery bound for the $\ell_{p,q}$ regularization problem in Theorem 2.1.

Lemma 2.4. Let $0 < q \leq 1 \leq p$, $\tau \geq 1$ and $x \in \mathbb{R}^n$, $\mathcal{N} := \mathcal{J}(x; N) \cup \mathcal{J}$ and $\mathcal{J}_k := \mathcal{J}_k(x; N)$ for $k = 1, \dots, R$. Then the following inequalities hold

$$\|x_{\mathcal{G}_{N^c}}\|_{p,\tau} \leq \sum_{k=1}^R \|x_{\mathcal{G}_{\mathcal{J}_k}}\|_{p,\tau} \leq N^{\frac{1}{\tau} - \frac{1}{q}} \|x_{\mathcal{G}_{\mathcal{J}^c}}\|_{p,q}.$$

Proof. By the definition of \mathcal{J}_k (cf. (2.14)), for all $j \in \mathcal{J}_k$, one has that

$$\|x_{\mathcal{G}_j}\|_p \leq \|x_{\mathcal{G}_i}\|_p, \quad \text{for } i \in \mathcal{J}_{k-1},$$

and thus

$$\|x_{\mathcal{G}_j}\|_p^q \leq \frac{1}{N} \sum_{i \in \mathcal{J}_{k-1}} \|x_{\mathcal{G}_i}\|_p^q = \frac{1}{N} \|x_{\mathcal{G}_{\mathcal{J}_{k-1}}}\|_{p,q}^q.$$

Consequently, we obtain that

$$\|x_{\mathcal{G}_{\mathcal{J}_k}}\|_{p,\tau}^\tau = \sum_{i \in \mathcal{J}_k} \|x_{\mathcal{G}_i}\|_p^\tau \leq N^{1-\tau/q} \|x_{\mathcal{G}_{\mathcal{J}_{k-1}}}\|_{p,q}^\tau.$$

Further by [32, Lemma 4.1] ($\tau \geq 1$ and $q \leq 1$), it follows that

$$\begin{aligned}\|x_{\mathcal{G}_{\mathcal{N}^c}}\|_{p,\tau} &= \left(\sum_{k=1}^R \sum_{i \in \mathcal{J}_k} \|x_{\mathcal{G}_i}\|_p^\tau \right)^{1/\tau} \leq \sum_{k=1}^R \|x_{\mathcal{G}_{\mathcal{J}_k}}\|_{p,\tau} \\ &\leq N^{\frac{1}{\tau} - \frac{1}{q}} \sum_{k=1}^R \|x_{\mathcal{G}_{\mathcal{J}_{k-1}}}\|_{p,q} \leq N^{\frac{1}{\tau} - \frac{1}{q}} \|x_{\mathcal{G}_{\mathcal{J}^c}}\|_{p,q}.\end{aligned}$$

The proof is complete. ■

2.3 Global recovery bound

In recent years, many articles have been devoted to establishing the oracle property and the global recovery bound for the ℓ_1 regularization problem (1.1) under the RIP or REC; see, e.g., [5, 9, 36, 51, 62]. However, to the best of our knowledge, few papers concentrate on investigating these properties for the lower-order regularization problem.

In the preceding subsections, we have introduced the general notion of (p, q) -GREC. Under the (p, q) -GREC(S, S), the solution of $Ax = b$ with group sparsity being S is unique. In this subsection, we will present the oracle property and the global recovery bound for the $\ell_{p,q}$ regularization problem (1.4) under the (p, q) -GREC. The oracle property provides an upper bound on the square error of the linear system and the violation of the true nonzero groups for each point in the level set of the objective function of (1.4)

$$\text{lev}_F(\bar{x}) := \{x \in \mathbb{R}^n : \|Ax - b\|_2^2 + \lambda \|x\|_{p,q}^q \leq \lambda \|\bar{x}\|_{p,q}^q\}.$$

Proposition 2.2. *Let $0 < q \leq 1 \leq p$, $S > 0$ and let the (p, q) -GREC(S, S) hold. Let \bar{x} be the unique solution of $Ax = b$ at a group sparsity level S , and \mathcal{S} be the index set of nonzero groups of \bar{x} . Let K be the smallest integer such that $2^{K-1}q \geq 1$. Then, for any $x^* \in \text{lev}_F(\bar{x})$, the following oracle inequality holds*

$$\|Ax^* - A\bar{x}\|_2^2 + \lambda \|x_{\mathcal{G}_{S^c}}^*\|_{p,q}^q \leq \lambda^{\frac{2}{2-q}} S^{(1-2^{-K})\frac{2}{2-q}} / \phi_{p,q}^{\frac{2q}{2-q}}(S, S). \quad (2.15)$$

Moreover, letting $\mathcal{N}_* := \mathcal{S} \cup \mathcal{S}(x^*; S)$, we have

$$\|x_{\mathcal{G}_{\mathcal{N}_*}}^* - \bar{x}_{\mathcal{G}_{\mathcal{N}_*}}\|_{p,2}^2 \leq \lambda^{\frac{2}{2-q}} S^{(1-2^{-K})\frac{2}{2-q}} / \phi_{p,q}^{\frac{4}{2-q}}(S, S).$$

Proof. Let $x^* \in \text{lev}_F(\bar{x})$. That is, $\|Ax^* - b\|_2^2 + \lambda \|x^*\|_{p,q}^q \leq \lambda \|\bar{x}\|_{p,q}^q$. By Lemmas 2.1(ii) and 2.2, one has that

$$\begin{aligned}\|Ax^* - A\bar{x}\|_2^2 + \lambda \|x_{\mathcal{G}_{S^c}}^*\|_{p,q}^q &\leq \lambda \|\bar{x}_{\mathcal{G}_S}\|_{p,q}^q - \lambda \|x_{\mathcal{G}_S}^*\|_{p,q}^q \\ &\leq \lambda \|\bar{x}_{\mathcal{G}_S} - x_{\mathcal{G}_S}^*\|_{p,q}^q \\ &\leq \lambda S^{1-2^{-K}} \|\bar{x}_{\mathcal{G}_S} - x_{\mathcal{G}_S}^*\|_{p,2}^q.\end{aligned} \quad (2.16)$$

Noting that

$$\|x_{\mathcal{G}_{S^c}}^* - \bar{x}_{\mathcal{G}_{S^c}}\|_{p,q}^q - \|x_{\mathcal{G}_S}^* - \bar{x}_{\mathcal{G}_S}\|_{p,q}^q \leq \|x_{\mathcal{G}_{S^c}}^*\|_{p,q}^q - (\|\bar{x}_{\mathcal{G}_S}\|_{p,q}^q - \|x_{\mathcal{G}_S}^*\|_{p,q}^q) = \|x^*\|_{p,q}^q - \|\bar{x}\|_{p,q}^q \leq 0.$$

Then (p, q) -GREC(S, S) implies that

$$\|\bar{x}_{\mathcal{G}_S} - x_{\mathcal{G}_S}^*\|_{p,2} \leq \|Ax^* - A\bar{x}\|_2 / \phi_{p,q}(S, S). \quad (2.17)$$

From (2.16) and (2.17), it follows that

$$\|Ax^* - A\bar{x}\|_2^2 + \lambda \|x_{\mathcal{G}_{Sc}}^*\|_{p,q}^q \leq \lambda S^{1-2^{-K}} \|Ax^* - A\bar{x}\|_2^q / \phi_{p,q}^q(S, S), \quad (2.18)$$

and consequently,

$$\|Ax^* - A\bar{x}\|_2 \leq \lambda^{\frac{1}{2-q}} S^{(1-2^{-K})/(2-q)} / \phi_{p,q}^{\frac{q}{2-q}}(S, S). \quad (2.19)$$

Therefore, by (2.18) and (2.19), we arrive at the oracle inequality (2.15). Furthermore, by the definition of \mathcal{N}_* , (p, q) -GREC(S, S) implies that

$$\|x_{\mathcal{G}_{N_*}}^* - \bar{x}_{\mathcal{G}_{N_*}}\|_{p,2}^2 \leq \|Ax^* - A\bar{x}\|_2^2 / \phi_{p,q}^2(S, S) \leq \lambda^{\frac{2}{2-q}} S^{(1-2^{-K})\frac{2}{2-q}} / \phi_{p,q}^{\frac{4}{2-q}}(S, S).$$

The proof is complete. ■

One of the main results of this section is presented as follows, where we establish the global recovery bound for the $\ell_{p,q}$ regularization problem under the (p, q) -GREC. We will apply oracle inequality (2.15) and Lemma 2.4 in our proof.

Theorem 2.1. *Let $0 < q \leq 1 \leq p \leq 2$, $S > 0$ and let the (p, q) -GREC(S, S) hold. Let \bar{x} be the unique solution of $Ax = b$ at a group sparsity level S , and \mathcal{S} be the index set of nonzero groups of \bar{x} . Let K be the smallest integer such that $2^{K-1}q \geq 1$. Then, for any $x^* \in \text{lev}_F(\bar{x})$, the following global recovery bound for (1.4) holds*

$$\|x^* - \bar{x}\|_2^2 \leq 2\lambda^{\frac{2}{2-q}} S^{\frac{q-2}{q} + (1-2^{-K})\frac{4}{q(2-q)}} / \phi_{p,q}^{\frac{4}{2-q}}(S, S). \quad (2.20)$$

More precisely,

$$\|x^* - \bar{x}\|_2^2 \leq \begin{cases} O\left(\lambda^{\frac{2}{2-q}} S\right), & 2^{K-1}q = 1, \\ O\left(\lambda^{\frac{2}{2-q}} S^{\frac{3-q}{2-q}}\right), & 2^{K-1}q > 1. \end{cases} \quad (2.21)$$

Proof. Let $\mathcal{N}_* := \mathcal{S} \cup \mathcal{S}(x^*; S)$ as in Proposition 2.2. Since $p \leq 2$, it follows from Lemma 2.4 and Proposition 2.2 that

$$\|x_{\mathcal{G}_{N_*}^c}^*\|_2^2 \leq \|x_{\mathcal{G}_{N_*}^c}^*\|_{p,2}^2 \leq S^{1-2/q} \|x_{\mathcal{G}_{Sc}}^*\|_{p,q}^2 \leq \lambda^{\frac{2}{2-q}} S^{\frac{q-2}{q} + (1-2^{-K})\frac{4}{q(2-q)}} / \phi_{p,q}^{\frac{4}{2-q}}(S, S).$$

Then by Proposition 2.2, one has that

$$\begin{aligned} \|x^* - \bar{x}\|_2^2 &= \|x_{\mathcal{G}_{N_*}}^* - \bar{x}_{\mathcal{G}_{N_*}}\|_2^2 + \|x_{\mathcal{G}_{N_*}^c}^*\|_2^2 \\ &\leq \lambda^{\frac{2}{2-q}} S^{(1-2^{-K})\frac{2}{2-q}} / \phi_{p,q}^{\frac{4}{2-q}}(S, S) + \lambda^{\frac{2}{2-q}} S^{\frac{q-2}{q} + (1-2^{-K})\frac{4}{q(2-q)}} / \phi_{p,q}^{\frac{4}{2-q}}(S, S) \\ &\leq 2\lambda^{\frac{2}{2-q}} S^{\frac{q-2}{q} + (1-2^{-K})\frac{4}{q(2-q)}} / \phi_{p,q}^{\frac{4}{2-q}}(S, S), \end{aligned}$$

where the last inequality follows from the fact that $2^{K-1}q \geq 1$. In particular, if $2^{K-1}q = 1$, then $\frac{q-2}{q} + (1 - 2^{-K}) \frac{4}{q(2-q)} = 1$ and thus

$$\|x^* - \bar{x}\|_2^2 \leq O\left(\lambda^{\frac{2}{2-q}} S\right).$$

If $2^{K-1}q > 1$, then $2^{K-2}q < 1$. Hence, $\frac{q-2}{q} + (1 - 2^{-K}) \frac{4}{q(2-q)} < \frac{3-q}{2-q}$, and consequently

$$\|x^* - \bar{x}\|_2^2 \leq O\left(\lambda^{\frac{2}{2-q}} S^{\frac{3-q}{2-q}}\right).$$

The proof is complete. ■

Theorem 2.1 is an important theoretical result in that it provides the global recovery bound (2.21) for the general $\ell_{p,q}$ regularization problem (1.4). In particular, when x^* is a global optimal solution of (1.4) as assumed in [51], Theorem 2.1 provides an upper bound on the distance from the optimal solution x^* to the true sparse solution \bar{x} . This result is significantly different from the previous works [5, 9, 36, 62], our result here is deterministic and does not involve any kind of randomization, and thus does not have a nonzero probability of failure. The following two corollaries show the recovery bounds for the group Lasso and $\ell_{p,1/2}$ regularization problem.

Corollary 2.1. *Let \bar{x} be a solution of $Ax = b$, and S be the group sparsity of \bar{x} . Let $1 \leq p \leq 2$, and let $x^*(\ell_{p,1})$ be an optimal solution of the $\ell_{p,1}$ regularization problem. Suppose that the $(p, 1)$ -GREC(S, S) holds. Then the following recovery bound for the $\ell_{p,1}$ regularization problem holds*

$$\|x^*(\ell_{p,1}) - \bar{x}\|_2^2 \leq O(\lambda^2 S). \quad (2.22)$$

Corollary 2.1 extends the study of the ℓ_1 regularization problem (Lasso) in [5, 51], where the recovery bound is given by

$$\|x^*(\ell_1) - \bar{x}\|_2^2 \leq O(\lambda^2 s). \quad (2.23)$$

Comparing with (2.23), Corollary 2.1 provides the theoretical evidence for the phenomenon that exploiting the group sparsity structure can enhance the recovery performance when $S \ll s$.

Corollary 2.2. *Let \bar{x} be a solution of $Ax = b$, and S be the group sparsity of \bar{x} . Let $1 \leq p \leq 2$, and let $x^*(\ell_{p,1/2})$ be a global optimal solution of the $\ell_{p,1/2}$ regularization problem. Suppose that the $(p, 1/2)$ -GREC(S, S) holds. Then the following global recovery bound for the $\ell_{p,1/2}$ regularization problem holds*

$$\|x^*(\ell_{p,1/2}) - \bar{x}\|_2^2 \leq O\left(\lambda^{4/3} S\right).$$

In particular, when $n_{max} = 1$, that is, the problem is in absence of the group structure, Corollary 2.2 exhibits the following global recovery bound for the $\ell_{1/2}$ regularization problem

$$\|x^*(\ell_{1/2}) - \bar{x}\|_2^2 \leq O\left(\lambda^{4/3}s\right). \quad (2.24)$$

Bound (2.22) seems to be the best one in a sense that all other bounds $O\left(\lambda^{\frac{2}{2-q}}S\right)$ are not better than it when $q < 1$ and the $(p, 1)$ -GREC(S, S) holds. However, when the $(p, 1)$ -GREC(S, S) does not hold but $(p, 1/2)$ -GREC(S, S) holds, we illustrate by an example that (2.22) or (2.23) does not hold but (2.24) does and is also tight. We will testify the recovery bound (2.24) by using a global optimization method.

Example 2.2. By letting $a = \tilde{a} = 2$ and $c = \tilde{c} = 1$ in Example 2.1, we consider the following matrix:

$$A = \begin{pmatrix} 2 & 3 & 1 \\ 2 & 1 & 3 \end{pmatrix}.$$

We assume $b = (2, 2)^\top$ and then a true solution of $Ax = b$ is $\bar{x} = (1, 0, 0)^\top$. Denoting $x := (x_1, x_2, x_3)^\top$, the objective function associated with the ℓ_1 regularization problem (1.1) is

$$\begin{aligned} F(x) &:= \|Ax - b\|_2^2 + \lambda\|x\|_1 \\ &= (2x_1 + 3x_2 + x_3 - 2)^2 + (2x_1 + x_2 + 3x_3 - 2)^2 + \lambda(|x_1| + |x_2| + |x_3|). \end{aligned}$$

Let $x^*(\ell_1) := (x_1^*, x_2^*, x_3^*)^\top$ be an optimal solution of problem (1.1). Without loss of generality, we assume $\lambda \leq 1$. The necessary condition of $x^*(\ell_1)$ being an optimal solution of (1.1) is $0 \in \partial F(x^*(\ell_1))$, that is,

$$0 \in 16x_1^* + 16x_2^* + 16x_3^* - 16 + \lambda\partial|x_1^*|, \quad (2.25a)$$

$$0 \in 16x_1^* + 20x_2^* + 12x_3^* - 16 + \lambda\partial|x_2^*|, \quad (2.25b)$$

$$0 \in 16x_1^* + 12x_2^* + 20x_3^* - 16 + \lambda\partial|x_3^*|, \quad (2.25c)$$

where $\partial|\mu| := \begin{cases} \text{sgn}(\mu), & \mu \neq 0, \\ [-1, 1], & \mu = 0. \end{cases}$

We first show that $x_i^* \geq 0$ for $i = 1, 2, 3$ by contradiction. Indeed, if $x_1^* < 0$, (2.25a) reduces to

$$16x_1^* + 16x_2^* + 16x_3^* - 16 = \lambda.$$

Summing (2.25b) and (2.25c), we further have

$$\lambda = 16x_1^* + 16x_2^* + 16x_3^* - 16 \in -\frac{\lambda}{2}(\partial|x_2^*| + \partial|x_3^*|),$$

which implies that $x_2^* \leq 0$ and $x_3^* \leq 0$. Hence, it follows that $F(x^*) > F(0)$, which indicates that x^* is not an optimal solution of (1.1), and thus, $x_1^* < 0$ is impossible. Similarly, we can show that $x_2^* \geq 0$ and $x_3^* \geq 0$.

Next, we find the optimal solution $x^*(\ell_1)$ by only considering $x^*(\ell_1) \geq 0$. It is easy to obtain that the solution of (2.25) and the corresponding objective value associated with (1.1) can be represented respectively by

$$x_1^* = 1 - \frac{\lambda}{16} - 2x_3^*, \quad x_2^* = x_3^* \left(0 \leq x_3^* \leq \frac{1}{2} - \frac{\lambda}{32} \right), \quad \text{and} \quad F(x^*(\ell_1)) = \lambda - \frac{\lambda^2}{32}.$$

Hence, $x^*(\ell_1) := (0, \frac{1}{2} - \frac{\lambda}{32}, \frac{1}{2} - \frac{\lambda}{32})^\top$ is an optimal solution of problem (1.1). The estimated error for such $x^*(\ell_1)$ is

$$\|x^*(\ell_1) - \bar{x}\|_2^2 = 1 + \frac{1}{2} \left(1 - \frac{\lambda}{16} \right)^2 > 1,$$

which does not meet the recovery bound (2.23) for each $\lambda \leq 1$.

It is revealed from Example 2.1 that this matrix A satisfies the $1/2\text{-REC}(1,1)$. Then the hypothesis of Corollary 2.2 is verified, and thus, Corollary 2.2 is applicable to establishing the recovery bound (2.24) for the $\ell_{1/2}$ regularization problem. Even though we cannot obtain the closed-form solution of this nonconvex $\ell_{1/2}$ regularization problem, as it is of only 3-dimensions, we use a global optimization method, the filled function method [28], to find the global optimal solution $x^*(\ell_{1/2})$ and thus to testify the recovery bound (2.24). This is done by computing the $\ell_{1/2}$ regularization problem for many λ to plot the curve $\|x^*(\ell_{1/2}) - \bar{x}\|_2^2$. Figure 1 illustrates the variation of the estimated error $\|x^*(\ell_{1/2}) - \bar{x}\|_2^2$ and the bound $2\lambda^{4/3}$ (that is the right-hand side of (2.20), where $S = 1$ and $\phi_{1/2}(1,1) \leq 1$ (cf. Example 2.1)), when varying the regularization parameter λ from 10^{-8} to 1.

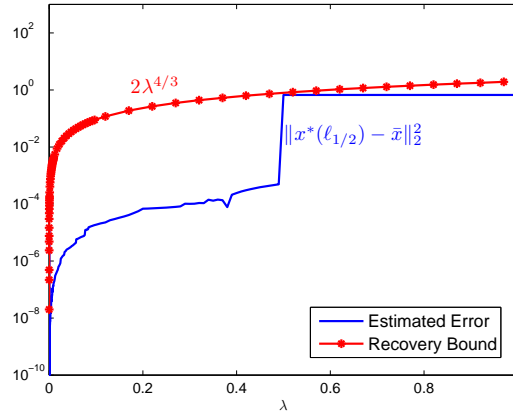


Figure 1: The illustration of the recovery bound (2.20) and estimated error.

2.4 Local recovery bound

In the preceding subsection, we provided the global analysis of the recovery bound for the $\ell_{p,q}$ regularization problem under the (p,q) -GREC; see Theorem 2.1. One can also see from

Figure 1 that the global recovery bound (2.24) is tight for the $\ell_{1/2}$ regularization problem as the curves come together at $\lambda \simeq 0.5$, but there is still a big gap for the improvement when λ is small.

This subsection is devoted to providing a local analysis of the recovery bound for the $\ell_{p,q}$ regularization problem by virtue of the technique of variational analysis [45]. For $x \in \mathbb{R}^n$ and $\delta \in \mathbb{R}_+$, $\mathbf{B}(x, \delta)$ denotes the open ball of radius δ centered at x . For a lower semi-continuous (lsc) function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $x, w \in \mathbb{R}^n$, the subderivative of f at x along the direction w is defined by

$$df(\bar{x})(w) := \liminf_{\tau \downarrow 0, w' \rightarrow w} \frac{f(\bar{x} + \tau w') - f(\bar{x})}{\tau}.$$

To begin with, we show in the following lemma a significant advantage of lower-order regularization over the ℓ_1 regularization: the lower-order regularization term can easily induce the sparsity of the local minimum.

Lemma 2.5. *Let $0 < q < 1 \leq p$. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a lsc function and $df(0)(0) = 0$. Then the function $F := f + \lambda \|\cdot\|_{p,q}^q$ has a local minimum at 0 with the first-order growth condition being fulfilled, i.e., there exist some $\epsilon > 0$ and $\delta > 0$ such that*

$$F(x) \geq F(0) + \epsilon \|x\|_2 \quad \forall x \in \mathbf{B}(0, \delta).$$

Proof. Let $\varphi := \lambda \|\cdot\|_{p,q}^q$ and then $F = f + \varphi$. Since φ is grouped separable, by [45, Proposition 10.5], it follows from the definition that $d\varphi(0) = \delta_{\{0\}}$ (where δ_X is the indicator function of X). Applying [45, Proposition 10.9], it follows that

$$dF(0) \geq df(0) + \delta_{\{0\}}. \quad (2.26)$$

Since f is finite and $df(0)(0) = 0$, its subderivative $df(0)$ is proper (cf. [45, Exercise 3.19]). Since further $df(0)(0) = 0$, it yields that $df(0) + \delta_{\{0\}} = \delta_{\{0\}}$. Thus by (2.26), we obtain that $dF(0) \geq \delta_{\{0\}}$. Therefore, by definition, there exist some $\epsilon > 0$ and $\delta > 0$ such that

$$F(x) \geq F(0) + \epsilon \|x\|_2 \quad \forall x \in \mathbf{B}(0, \delta).$$

The proof is complete. ■

With the help of the above lemma, we can present in the following a local version of the recovery bound. This is done by constructing a path of local minima depending on the regularization parameter λ for the regularization problem, which starts from a sparse solution of the original problem and shares the same support as this sparse solution has, resulting in a sharper bound in terms of λ^2 .

Theorem 2.2. *Let \bar{x} be a solution of $Ax = b$, S be the group sparsity of \bar{x} , and B be a submatrix of A consisting of its columns corresponding to the active components of \bar{x} . Suppose that any nonzero group of \bar{x} is active, and that the columns of A corresponding to*

the active components of \bar{x} are linearly independent. Let $0 < q < 1 \leq p$. Then there exist $\kappa > 0$ and a path of local minima of problem (1.4), $x^*(\lambda)$, such that

$$\|x^*(\lambda) - \bar{x}\|_2^2 \leq \lambda^2 q^2 S \|(B^\top B)^{-1}\|^2 \max_{\bar{x}_{\mathcal{G}_i} \neq 0} \left(\|\bar{x}_{\mathcal{G}_i}\|_p^{2(q-p)} \|\bar{x}_{\mathcal{G}_i}\|_{2p-2}^{2p-2} \right) \quad \forall \lambda < \kappa.$$

Proof. Without loss of generality, we let \bar{x} be of structure $\bar{x} = (\bar{z}^\top, 0)^\top$ with

$$\bar{z} = (\bar{x}_{\mathcal{G}_1}^\top, \dots, \bar{x}_{\mathcal{G}_S}^\top)^\top \text{ and } \bar{x}_{\mathcal{G}_i} \neq \mathbf{0} \text{ for } i = 1, \dots, S.$$

and let s be the sparsity of \bar{x} . Let $A = (B, D)$ with B being the submatrix involving the first s columns of A (corresponding to the active components of \bar{x}). By the assumption, we have that B is of full column rank and thus $B^\top B$ is invertible. In this setting, the linear relation $A\bar{x} = b$ reduces to $B\bar{z} = b$. The proof of this theorem is divided into the three steps:

- (a) construct a smooth path from \bar{x} by the implicit function theorem;
- (b) validate that every point of the constructed path is a local minimum of (1.4); and
- (c) establish the recovery bound for the constructed path.

First, to show (a), we define $H : \mathbb{R}^{s+1} \rightarrow \mathbb{R}^s$ by

$$H(z, \lambda) = 2B^\top (Bz - b) + \lambda q \begin{pmatrix} \|z_{\mathcal{G}_1}\|_p^{q-p} \sigma(z_{\mathcal{G}_1}) \\ \vdots \\ \|z_{\mathcal{G}_S}\|_p^{q-p} \sigma(z_{\mathcal{G}_S}) \end{pmatrix},$$

where $\sigma(z_{\mathcal{G}_i}) = \text{vector}(|z_j|^{p-1} \text{sign}(z_j))_{\mathcal{G}_i}$, denoting a vector consisting of $|z_j|^{p-1} \text{sign}(z_j)$ for all $j \in \mathcal{G}_i$. Let $\bar{\delta} > 0$ be sufficiently small such that $\text{sign}(z) = \text{sign}(\bar{z})$ for each $z \in \mathbf{B}(\bar{z}, \bar{\delta})$ and thus H is smooth on $\mathbf{B}(\bar{z}, \bar{\delta}) \times \mathbb{R}$. Note that $H(\bar{z}, 0) = 0$ and $\frac{\partial H}{\partial z}(\bar{z}, 0) = 2B^\top B$. By the implicit function theorem [46], there exist some $\kappa > 0$, $\delta \in (0, \bar{\delta})$ and a unique smooth function $\xi : (-\kappa, \kappa) \rightarrow \mathbf{B}(\bar{z}, \delta)$ such that

$$\{(z, \lambda) \in \mathbf{B}(\bar{z}, \bar{\delta}) \times (-\kappa, \kappa) : H(z, \lambda) = 0\} = \{(\xi(\lambda), \lambda) : \lambda \in (-\kappa, \kappa)\}, \quad (2.27)$$

and

$$\frac{d\xi}{d\lambda} = -q \left(2B^\top B + \lambda q \begin{pmatrix} M_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & M_S \end{pmatrix} \right)^{-1} \begin{pmatrix} \|\xi(\lambda)_{\mathcal{G}_1}\|_p^{q-p} \sigma(\xi(\lambda)_{\mathcal{G}_1}) \\ \vdots \\ \|\xi(\lambda)_{\mathcal{G}_S}\|_p^{q-p} \sigma(\xi(\lambda)_{\mathcal{G}_S}) \end{pmatrix}, \quad (2.28)$$

where M_i for each $i = 1, \dots, S$ is denoted by

$$M_i = (q - p) \|\xi(\lambda)_{\mathcal{G}_i}\|_p^{q-2p} (\sigma(\xi(\lambda)_{\mathcal{G}_i})) (\sigma(\xi(\lambda)_{\mathcal{G}_i}))^\top + (p - 1) \|\xi(\lambda)_{\mathcal{G}_i}\|_p^{q-p} \text{diag}(|\xi(\lambda)_j|^{p-2}),$$

and $\text{diag}(|\xi(\lambda)_j|^{p-2})$ is a diagonal matrix generated by vector $(|\xi(\lambda)_j|^{p-2})$. Thus, due to (2.27) and (2.28), we have constructed a smooth path $\xi(\lambda)$ near \bar{z} , $\lambda \in (-\kappa, \kappa)$, be such that

$$2B^\top(B\xi(\lambda) - b) + \lambda q \begin{pmatrix} \|\xi(\lambda)_{\mathcal{G}_1}\|_p^{q-p} \sigma(\xi(\lambda)_{\mathcal{G}_1}) \\ \vdots \\ \|\xi(\lambda)_{\mathcal{G}_S}\|_p^{q-p} \sigma(\xi(\lambda)_{\mathcal{G}_S}) \end{pmatrix} = 0, \quad (2.29)$$

and that

$$2B^\top B + \lambda q \begin{pmatrix} M_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & M_S \end{pmatrix} \succ 0. \quad (2.30)$$

For fixed $\lambda \in (-\kappa, \kappa)$, let $x^*(\lambda) := (\xi(\lambda)^\top, 0)^\top$. To verify (b), we prove that such $x^*(\lambda)$, with $\xi(\lambda)$ satisfying (2.29) and (2.30), is a local minimum of problem (1.4). Let $h(z) := \|Bz - b\|_2^2 + \lambda \|z\|_{p,q}^q$. Note that $h(\xi(\lambda)) = \|Ax^*(\lambda) - b\|_2^2 + \lambda \|x^*(\lambda)\|_{p,q}^q$ and that h is smooth around $\xi(\lambda)$. By noting that $\xi(\lambda)$ satisfies (2.29) and (2.30) (the first- and second- derivative of h at $\xi(\lambda)$), one has that h satisfies the second-order growth condition at $\xi(\lambda)$, that is, there exist some $\epsilon_\lambda > 0$ and $\delta_\lambda > 0$ such that

$$h(z) \geq h(\xi(\lambda)) + 2\epsilon_\lambda \|z - \xi(\lambda)\|_2^2 \quad \forall z \in \mathbf{B}(\xi(\lambda), \delta_\lambda). \quad (2.31)$$

In what follows, let $\epsilon_\lambda > 0$ and $\delta_\lambda > 0$ be given as above, and select $\epsilon_0 > 0$ such that

$$\sqrt{\epsilon_\lambda \epsilon_0} - \|B\| \|D\| > 0. \quad (2.32)$$

According to Lemma 2.5 (with $\|D \cdot\|_2^2 + 2\langle B\xi(\lambda) - b, D \cdot \rangle - 2\epsilon_0 \|\cdot\|_2^2$ in place of f), there exists some $\delta_0 > 0$ such that

$$\|Dy\|_2^2 + 2\langle B\xi(\lambda) - b, Dy \rangle - 2\epsilon_0 \|y\|_2^2 + \lambda \|y\|_{p,q}^p \geq 0 \quad \forall y \in \mathbf{B}(0, \delta_0). \quad (2.33)$$

Thus, for each $x := (z, y) \in \mathbf{B}(\xi(\lambda), \delta_\lambda) \times \mathbf{B}(0, \delta_0)$, it follows that

$$\begin{aligned} \|Ax - b\|_2^2 + \lambda \|x\|_{p,q}^q &= \|Bz - b + Dy\|_2^2 + \lambda \|z\|_{p,q}^q + \lambda \|y\|_{p,q}^q \\ &= \|Bz - b\|_2^2 + \lambda \|z\|_{p,q}^q + \|Dy\|_2^2 + 2\langle Bz - b, Dy \rangle + \lambda \|y\|_{p,q}^q \\ &= h(z) + \|Dy\|_2^2 + 2\langle B\xi(\lambda) - b, Dy \rangle + \lambda \|y\|_{p,q}^q + 2\langle B(z - \xi(\lambda)), Dy \rangle. \end{aligned}$$

By (2.31) and (2.33), it yields that

$$\begin{aligned} \|Ax - b\|_2^2 + \lambda \|x\|_{p,q}^q &\geq h(\xi(\lambda)) + 2\epsilon_\lambda \|z - \xi(\lambda)\|_2^2 + 2\epsilon_0 \|y\|_2^2 + 2\langle B(z - \xi(\lambda)), Dy \rangle \\ &\geq h(\xi(\lambda)) + 4\sqrt{\epsilon_\lambda \epsilon_0} \|z - \xi(\lambda)\|_2 \|y\|_2 - 2\|B\| \|D\| \|z - \xi(\lambda)\|_2 \|y\|_2 \\ &= \|Ax^*(\lambda) - b\|_2^2 + \lambda \|x^*(\lambda)\|_{p,q}^q + 2(2\sqrt{\epsilon_\lambda \epsilon_0} - \|B\| \|D\|) \|z - \xi(\lambda)\|_2 \|y\|_2 \\ &\geq \|Ax^*(\lambda) - b\|_2^2 + \lambda \|x^*(\lambda)\|_{p,q}^q, \end{aligned}$$

where the last inequality follows from (2.32). Hence $x^*(\lambda)$ is a local minimum of problem (1.4), and (b) is verified.

Finally, we check (c) by providing an upper bound on the distance from $\xi(\lambda)$ to \bar{z} . By (2.29), one has that

$$\xi(\lambda) - \bar{z} = -\frac{\lambda q}{2}((B^\top B)^{-1}) \begin{pmatrix} \|\xi(\lambda)_{\mathcal{G}_1}\|_p^{q-p} \sigma(\xi(\lambda)_{\mathcal{G}_1}) \\ \vdots \\ \|\xi(\lambda)_{\mathcal{G}_S}\|_p^{q-p} \sigma(\xi(\lambda)_{\mathcal{G}_S}) \end{pmatrix}. \quad (2.34)$$

As $\{\xi(\lambda) : \lambda \in (-\kappa, \kappa)\} \subseteq \mathbf{B}(\bar{z}, \bar{\delta})$, without loss of generality, we assume for each $\lambda < \kappa$ that

$$\|\xi(\lambda)_{\mathcal{G}_i}\|_p^{2(q-p)} \leq 2\|\bar{z}_{\mathcal{G}_i}\|_p^{2(q-p)} \quad \text{and} \quad \|\xi(\lambda)_{\mathcal{G}_i}\|_{2p-2}^{2p-2} \leq 2\|\bar{z}_{\mathcal{G}_i}\|_{2p-2}^{2p-2} \quad \forall i = 1, \dots, S$$

(otherwise, we can choose a smaller $\bar{\delta}$). Recall $\sigma(\xi(\lambda)_{\mathcal{G}_i}) = \text{vector}(|\xi(\lambda)_j|^{p-1} \text{sign}(\xi(\lambda)_j))_{\mathcal{G}_i}$. We obtain from (2.34) that

$$\begin{aligned} \|\xi(\lambda) - \bar{z}\|_2^2 &\leq \frac{\lambda^2 q^2}{4} \|(B^\top B)^{-1}\|^2 \sum_{i=1}^S \left(\|\xi(\lambda)_{\mathcal{G}_i}\|_p^{2(q-p)} \sum_{j \in \mathcal{G}_i} |\xi(\lambda)_j|^{2p-2} \right) \\ &= \frac{\lambda^2 q^2}{4} \|(B^\top B)^{-1}\|^2 \sum_{i=1}^S \left(\|\xi(\lambda)_{\mathcal{G}_i}\|_p^{2(q-p)} \|\xi(\lambda)_{\mathcal{G}_i}\|_{2p-2}^{2p-2} \right) \\ &\leq \frac{\lambda^2 q^2}{4} \|(B^\top B)^{-1}\|^2 S \max_{i=1, \dots, S} \left(\|\xi(\lambda)_{\mathcal{G}_i}\|_p^{2(q-p)} \|\xi(\lambda)_{\mathcal{G}_i}\|_{2p-2}^{2p-2} \right) \\ &\leq \lambda^2 q^2 S \|(B^\top B)^{-1}\|^2 \max_{i=1, \dots, S} \left(\|\bar{z}_{\mathcal{G}_i}\|_p^{2(q-p)} \|\bar{z}_{\mathcal{G}_i}\|_{2p-2}^{2p-2} \right). \end{aligned}$$

Hence we arrive at that for each $\lambda < \kappa$

$$\|x^*(\lambda) - \bar{x}\|_2^2 = \|\xi(\lambda) - \bar{z}\|_2^2 \leq \lambda^2 q^2 S \|(B^\top B)^{-1}\|^2 \max_{\bar{x}_{\mathcal{G}_i} \neq 0} \left(\|\bar{x}_{\mathcal{G}_i}\|_p^{2(q-p)} \|\bar{x}_{\mathcal{G}_i}\|_{2p-2}^{2p-2} \right),$$

and the proof is complete. ■

Theorem 2.2 is a significant result in that it provides the uniform local recovery bound for all the $\ell_{p,q}$ regularization problems ($0 < q < 1$), which is

$$\|x_{p,q}^*(\lambda) - \bar{x}\|_2^2 \leq O(\lambda^2 S),$$

where $x_{p,q}^*(\lambda)$ is a local optimal solution of (1.4) (near \bar{x}). This bound improves the global recovery bound given in Theorem 2.1 (of order $\lambda^{\frac{2}{2-q}}$) and shares the same one with the $\ell_{p,1}$ regularization problem (group Lasso); see Corollary 2.1. It is worth noting that when $q = 1$ our proof technique is not working as Lemma 2.5 is false in this case.

3 Proximal gradient method for group sparse optimization

Many efficient algorithms have been proposed to solve the sparse optimization problem, and one of the most popular and practical algorithms is the proximal gradient method; see [3, 16, 39, 55] and references therein. It was reported in [3, 16, 39] that the proximal gradient method for solving the ℓ_1 regularization problem reduces to the well-known iterative soft thresholding algorithm and that the accelerated proximal gradient methods proposed in

[3, 39] have the convergence rate $O(1/k^2)$. Recently, the convergence theory of the proximal gradient method for solving the nonconvex regularization problem was studied under the framework of the Kurdyka-Łojasiewicz theory [1, 7], the majorization-minimization (MM) scheme [34], the coordinate gradient descent method [49] and the successive upper-bound minimization approach [44].

In this section, we apply the proximal gradient method to solve the group sparse optimization (PGM-GSO) via $\ell_{p,q}$ ($p \geq 1, q \geq 0$) regularization (1.4), which is stated as follows.

PGM-GSO

Select a stepsize v , start with an initial point $x_0 \in \mathbb{R}^n$, and generate a sequence $\{x^k\} \subseteq \mathbb{R}^n$ via the iteration

$$z^k = x^k - 2vA^\top(Ax^k - b), \quad (3.1)$$

$$x^{k+1} \in \operatorname{Arg} \min_{x \in \mathbb{R}^n} \left\{ \lambda \|x\|_{p,q}^q + \frac{1}{2v} \|x - z^k\|_2^2 \right\}. \quad (3.2)$$

We will obtain the analytical solutions of (3.2) for some specific p and q , and the linear convergence rate of the PGM-GSO. The convergence theory of the PGM-GSO falls in the framework of the Kurdyka-Łojasiewicz theory; see [1, 7]. In particular, following from [7, Theorem 1 and Proposition 3], the sequence generated by the PGM-GSO converges to a critical point, especially a global minimum when $q \geq 1$ and a local minimum when $q = 0$ (inspired by the idea in [6]), as summarized as follows.

Theorem 3.1. *Let $p \geq 1$. Suppose that the sequence $\{x^k\}$ is generated by the PGM-GSO with $v < \frac{1}{2}\|A\|_2^{-2}$. Then the following statements hold:*

- (i) *if $q \geq 1$, then $\{x^k\}$ converges to a global minimum of problem (1.4),*
- (ii) *if $q = 0$, then $\{x^k\}$ converges to a local minimum of problem (1.4), and*
- (iii) *if $0 < q < 1$, then $\{x^k\}$ converges to a critical point¹ of problem (1.4).*

3.1 Analytical solutions of (3.2)

Since the main computation of the proximal gradient method is the proximal step (3.2), it is significant to investigate the solutions of (3.2) for the specific applications. Note that $\|x\|_{p,q}^q$ and $\|x - z^k\|_2^2$ are both grouped separable. Then the proximal step (3.2) can be achieved parallelly in each group, and is equivalent to solve a cycle of low dimensional proximal optimization subproblems

$$x_{\mathcal{G}_i}^{k+1} \in \operatorname{Arg} \min_{x \in \mathbb{R}^{n_i}} \left\{ \lambda \|x_{\mathcal{G}_i}\|_p^q + \frac{1}{2v} \|x_{\mathcal{G}_i} - z_{\mathcal{G}_i}^k\|_2^2 \right\}, \text{ for } i = 1, \dots, r. \quad (3.3)$$

When p and q are given as some specific numbers, such as $p = 1, 2$ and $q = 0, 1/2, 2/3, 1$, the solution of subproblem (3.3) of each group can be given explicitly by an analytical formula, as shown in the following proposition.

¹A point x is said to be a critical point of F if 0 belongs to its limiting subdifferential at x [37].

Proposition 3.1. Let $z \in \mathbb{R}^l$, $v > 0$ and the proximal regularization be $Q_{p,q}(x) := \lambda \|x\|_p^q + \frac{1}{2v} \|x - z\|_2^2$. Then the proximal operator

$$P_{p,q}(z) \in \text{Arg min}_{x \in \mathbb{R}^l} \{Q_{p,q}(x)\}$$

has the following analytical formula:

(i) if $p = 2$ and $q = 1$, then

$$P_{2,1}(z) = \begin{cases} \left(1 - \frac{v\lambda}{\|z\|_2}\right) z, & \|z\|_2 > v\lambda, \\ 0, & \text{otherwise,} \end{cases} \quad (3.4)$$

(ii) if $p = 2$ and $q = 0$, then

$$P_{p,0}(z) = \begin{cases} z, & \|z\|_2 > \sqrt{2v\lambda}, \\ 0 \text{ or } z, & \|z\|_2 = \sqrt{2v\lambda}, \\ 0, & \|z\|_2 < \sqrt{2v\lambda}, \end{cases} \quad (3.5)$$

(iii) if $p = 2$ and $q = 1/2$, then

$$P_{2,1/2}(z) = \begin{cases} \frac{16\|z\|_2^{3/2} \cos^3\left(\frac{\pi}{3} - \frac{\psi(z)}{3}\right)}{3\sqrt{3}v\lambda + 16\|z\|_2^{3/2} \cos^3\left(\frac{\pi}{3} - \frac{\psi(z)}{3}\right)} z, & \|z\|_2 > \frac{3}{2}(v\lambda)^{2/3}, \\ 0 \text{ or } \frac{16\|z\|_2^{3/2} \cos^3\left(\frac{\pi}{3} - \frac{\psi(z)}{3}\right)}{3\sqrt{3}v\lambda + 16\|z\|_2^{3/2} \cos^3\left(\frac{\pi}{3} - \frac{\psi(z)}{3}\right)} z, & \|z\|_2 = \frac{3}{2}(v\lambda)^{2/3}, \\ 0, & \|z\|_2 < \frac{3}{2}(v\lambda)^{2/3}, \end{cases} \quad (3.6)$$

with

$$\psi(z) = \arccos\left(\frac{v\lambda}{4} \left(\frac{3}{\|z\|_2}\right)^{3/2}\right), \quad (3.7)$$

(iv) if $p = 1$ and $q = 1/2$, then

$$P_{1,1/2}(z) = \begin{cases} \tilde{z}, & Q_{1,1/2}(\tilde{z}) < Q_{1,1/2}(0), \\ 0 \text{ or } \tilde{z}, & Q_{1,1/2}(\tilde{z}) = Q_{1,1/2}(0), \\ 0, & Q_{1,1/2}(\tilde{z}) > Q_{1,1/2}(0), \end{cases} \quad (3.8)$$

with

$$\tilde{z} = z - \frac{\sqrt{3}v\lambda}{4\sqrt{\|z\|_1} \cos\left(\frac{\pi}{3} - \frac{\xi(z)}{3}\right)} \text{sign}(z), \quad \xi(z) = \arccos\left(\frac{v\lambda}{4} \left(\frac{3}{\|z\|_1}\right)^{3/2}\right),$$

(v) if $p = 2$ and $q = 2/3$, then

$$P_{2,2/3}(z) = \begin{cases} \frac{3(a^{3/2} + \sqrt{2\|z\|_2 - a^3})}{32v\lambda a^2 + 3(a^{3/2} + \sqrt{2\|z\|_2 - a^3})} z, & \|z\|_2 > 2\left(\frac{2}{3}v\lambda\right)^{3/4}, \\ 0 \text{ or } \frac{3(a^{3/2} + \sqrt{2\|z\|_2 - a^3})}{32v\lambda a^2 + 3(a^{3/2} + \sqrt{2\|z\|_2 - a^3})} z, & \|z\|_2 = 2\left(\frac{2}{3}v\lambda\right)^{3/4}, \\ 0, & \|z\|_2 < 2\left(\frac{2}{3}v\lambda\right)^{3/4}, \end{cases} \quad (3.9)$$

with

$$a = \frac{2}{\sqrt{3}}(2v\lambda)^{1/4} \left(\cosh \left(\frac{\varphi(z)}{3} \right) \right)^{1/2}, \quad \varphi(z) = \operatorname{arccosh} \left(\frac{27\|z\|_2^2}{16(2v\lambda)^{3/2}} \right), \quad (3.10)$$

(vi) if $p = 1$ and $q = 2/3$, then

$$P_{1,2/3}(z) = \begin{cases} \bar{z}, & Q_{1,2/3}(\bar{z}) < Q_{1,2/3}(0), \\ 0 \text{ or } \bar{z}, & Q_{1,2/3}(\bar{z}) = Q_{1,2/3}(0), \\ 0, & Q_{1,2/3}(\bar{z}) > Q_{1,2/3}(0), \end{cases} \quad (3.11)$$

with

$$\bar{z} = z - \frac{4v\lambda\bar{a}^{1/2}}{3(\bar{a}^{3/2} + \sqrt{2}\|z\|_1 - \bar{a}^3)} \operatorname{sign}(z),$$

and

$$\bar{a} = \frac{2}{\sqrt{3}}(2v\lambda l)^{1/4} \left(\cosh \left(\frac{\zeta(z)}{3} \right) \right)^{1/2}, \quad \zeta(z) = \operatorname{arccosh} \left(\frac{27\|z\|_1^2}{16(2v\lambda l)^{3/2}} \right).$$

Proof. Since the proximal regularization $Q_{p,q}(x) := \lambda\|x\|_p^q + \frac{1}{2v}\|x - z\|_2^2$ is non-differentiable only at 0, $P_{p,q}(z)$ must be 0 or some point $\tilde{x} (\neq 0)$ satisfying the first-order condition

$$\lambda q \|\tilde{x}\|_p^{q-p} \begin{pmatrix} |\tilde{x}_1|^{p-1} \operatorname{sign}(\tilde{x}_1) \\ \vdots \\ |\tilde{x}_l|^{p-1} \operatorname{sign}(\tilde{x}_l) \end{pmatrix} + \frac{1}{v}(\tilde{x} - z) = 0. \quad (3.12)$$

Thus, to derive the analytical formula of the proximal operator $P_{p,q}(z)$, we just need to calculate such \tilde{x} via (3.12), and then compare the objective function values $Q_{p,q}(\tilde{x})$ and $Q_{p,q}(0)$ to obtain the solution inducing the smaller value. The proofs of the six statements follow in the above routine, and we only provide the detailed proofs of (iii) and (v) as samples.

(iii) When $p = 2$ and $q = 1/2$, (3.12) reduces to

$$\frac{\lambda\tilde{x}}{2\|\tilde{x}\|_2^{3/2}} + \frac{1}{v}(\tilde{x} - z) = 0, \quad (3.13)$$

and consequently,

$$\|\tilde{x}\|_2^{3/2} - \|z\|_2\|\tilde{x}\|_2^{1/2} + \frac{1}{2}v\lambda = 0. \quad (3.14)$$

Denote $\eta = \|\tilde{x}\|_2^{1/2} > 0$. The equation (3.14) can be transformed into the following cubic algebraic equation

$$\eta^3 - \|z\|_2\eta + \frac{1}{2}v\lambda = 0. \quad (3.15)$$

Due to the hyperbolic solution of the cubic equation (see [47]), by denoting

$$r = 2\sqrt{\frac{\|z\|_2}{3}}, \quad \alpha = \arccos \left(\frac{v\lambda}{4} \left(\frac{3}{\|z\|_2} \right)^{3/2} \right) \quad (:= \psi(z)) \quad \text{and} \quad \beta = \operatorname{arccosh} \left(-\frac{v\lambda}{4} \left(\frac{3}{\|z\|_2} \right)^{3/2} \right),$$

the solution of (3.15) can be expressed as the follows.

(1) If $0 \leq \|z\|_2 \leq 3 \left(\frac{v\lambda}{4}\right)^{2/3}$, then the three roots of (3.15) are given by

$$\eta_1 = r \cosh \frac{\beta}{3}, \quad \eta_2 = -\frac{r}{2} \cosh \frac{\beta}{3} + i \frac{\sqrt{3}r}{2} \sinh \frac{\beta}{3}, \quad \eta_3 = -\frac{r}{2} \cosh \frac{\beta}{3} - i \frac{\sqrt{3}r}{2} \sinh \frac{\beta}{3},$$

where i denotes the imaginary unit. However, this β does not exist since the value of hyperbolic cosine must be positive. Thus, in this case, $P_{2,1/2}(z) = 0$.

(2) If $\|z\|_2 > 3 \left(\frac{v\lambda}{4}\right)^{2/3}$, then the three roots of (3.15) are

$$\eta_1 = r \cos \left(\frac{\pi}{3} - \frac{\alpha}{3} \right), \quad \eta_2 = -r \sin \left(\frac{\pi}{2} - \frac{\alpha}{3} \right), \quad \eta_3 = -r \cos \left(\frac{2\pi}{3} - \frac{\alpha}{3} \right).$$

The unique positive solution of (3.15) is $\|\tilde{x}\|_2^{1/2} = \eta_1$, and thus, the unique solution of (3.13) is given by

$$\tilde{x} = \frac{2\eta_1^3}{v\lambda + 2\eta_1^3} z = \frac{16\|z\|_2^{3/2} \cos^3 \left(\frac{\pi}{3} - \frac{\psi(z)}{3} \right)}{3\sqrt{3}v\lambda + 16\|z\|_2^{3/2} \cos^3 \left(\frac{\pi}{3} - \frac{\psi(z)}{3} \right)} z.$$

Finally, we compare the objective function values $Q_{2,1/2}(\tilde{x})$ and $Q_{2,1/2}(0)$. For this purpose, when $\|z\|_2 > 3 \left(\frac{v\lambda}{4}\right)^{2/3}$, we define

$$\begin{aligned} H(\|z\|_2) &:= \frac{v}{\|\tilde{x}\|_2} (Q_{2,1/2}(0) - Q_{2,1/2}(\tilde{x})) \\ &= \frac{v}{\|\tilde{x}\|_2} \left(\frac{1}{2v} \|z\|_2^2 - \lambda \|\tilde{x}\|_2^{1/2} - \frac{1}{2v} \|\tilde{x} - z\|_2^2 \right) \\ &= \|z\|_2 - \frac{\|\tilde{x}\|_2^2 + 2v\lambda \|\tilde{x}\|_2^{1/2}}{2\|\tilde{x}\|_2} \\ &= \frac{1}{2} \|z\|_2 - \frac{3}{4} v\lambda \|\tilde{x}\|_2^{-1/2}, \end{aligned}$$

where the third equality holds since that \tilde{x} is proportional to z , and fourth equality follows from (3.14). Since both $\|z\|_2$ and $\|\tilde{x}\|_2$ are strictly increasing on $\|z\|_2$, $H(\|z\|_2)$ is also strictly increasing when $\|z\|_2 > 3 \left(\frac{v\lambda}{4}\right)^{2/3}$. Thus the unique solution of $H(\|z\|_2) = 0$ satisfies

$$\|z\|_2 \|\tilde{x}\|_2^{1/2} = \frac{3}{2} v\lambda,$$

and further, (3.14) implies that the solution of $H(\|z\|_2) = 0$ is

$$\|z\|_2 = \frac{3}{2} (v\lambda)^{2/3}.$$

Therefore, we arrive at the formulae (3.6) and (3.7).

(v) When $p = 2$ and $q = 2/3$, (3.12) reduces to

$$\frac{2\lambda\tilde{x}}{3\|\tilde{x}\|_2^{4/3}} + \frac{1}{v}(\tilde{x} - z) = 0, \quad (3.16)$$

and consequently,

$$\|\tilde{x}\|_2^{4/3} - \|z\|_2 \|\tilde{x}\|_2^{1/3} + \frac{2}{3}v\lambda = 0. \quad (3.17)$$

Denote $\eta = \|\tilde{x}\|_2^{1/3} > 0$ and $h(t) = t^4 - \|z\|_2 t + \frac{2}{3}v\lambda$. Thus, η is the positive solution of $h(t) = 0$. Next, we seek η by the method of undetermined coefficients. Assume that

$$h(t) = t^4 - \|z\|_2 t + \frac{2}{3}v\lambda = (t^2 + at + b)(t^2 + ct + d), \quad \text{where } a, b, c, d \in \mathbb{R}. \quad (3.18)$$

By expansion and comparison, we have

$$a + c = 0, \quad b + d + ac = 0, \quad ad + bc = -\|z\|_2, \quad bd = \frac{2}{3}v\lambda,$$

and thus,

$$c = -a, b = \frac{1}{2} \left(a^2 + \frac{\|z\|_2}{a} \right), d = \frac{1}{2} \left(a^2 - \frac{\|z\|_2}{a} \right), bd = \frac{1}{4} \left(a^4 - \frac{\|z\|_2^2}{a^2} \right) = \frac{2}{3}v\lambda. \quad (3.19)$$

By letting $M = a^2$, the last one of the above equalities reduces to the following cubic algebraic equation

$$M^3 - \frac{8}{3}v\lambda M - \|z\|_2^2 = 0. \quad (3.20)$$

According to the Cardano formula for the cubic equation, the root of (3.20) can be represented by

$$a^2 = M = \left(\frac{\|z\|_2^2}{2} + \sqrt{\frac{\|z\|_2^4}{4} - \left(\frac{8}{9}v\lambda \right)^3} \right)^{1/3} + \left(\frac{\|z\|_2^2}{2} - \sqrt{\frac{\|z\|_2^4}{4} - \left(\frac{8}{9}v\lambda \right)^3} \right)^{1/3},$$

which can also be reformulated in the following hyperbolic form (see [47])

$$a^2 = M = \frac{4}{3}\sqrt{2v\lambda} \cosh \left(\frac{\varphi(z)}{3} \right), \quad (3.21)$$

where $\varphi(z)$ is given by (3.10). By (3.18) and (3.19), we have that η , the positive root of $h(t) = 0$, satisfies

$$\eta^2 + a\eta + \frac{1}{2} \left(a^2 + \frac{\|z\|_2}{a} \right) = 0 \quad \text{or} \quad \eta^2 - a\eta + \frac{1}{2} \left(a^2 - \frac{\|z\|_2}{a} \right) = 0.$$

Hence, the real roots of the above equations, that is, the real roots of $h(t) = 0$, are

$$\eta_1 = \frac{1}{2} \left(|a| + \sqrt{\frac{2\|z\|_2}{|a|} - a^2} \right), \quad \eta_2 = \frac{1}{2} \left(|a| - \sqrt{\frac{2\|z\|_2}{|a|} - a^2} \right). \quad (3.22)$$

It is easy to see that $\eta_1 > \eta_2$ and that η_2 should be discarded as it induces the saddle point rather than a minimum (since $h(t) > 0$ when $t < \eta_2$). Thus, by (3.16), (3.21) and (3.22), one has

$$\tilde{x} = \frac{3\eta_1^4}{2v\lambda + 3\eta_1^4} z = \frac{3 \left(a^{3/2} + \sqrt{2\|z\|_2 - a^3} \right)}{32v\lambda a^2 + 3 \left(a^{3/2} + \sqrt{2\|z\|_2 - a^3} \right)} z,$$

where a is given by (3.10). Finally, we compare the objective function values $Q_{2,2/3}(\tilde{x})$ and $Q_{2,2/3}(0)$. For this purpose, we define

$$\begin{aligned} H(\|z\|_2) &:= \frac{v}{\|\tilde{x}\|_2} (Q_{2,2/3}(0) - Q_{2,2/3}(\tilde{x})) \\ &= \frac{v}{\|\tilde{x}\|_2} \left(\frac{1}{2v} \|z\|_2^2 - \lambda \|\tilde{x}\|_2^{2/3} - \frac{1}{2v} \|\tilde{x} - z\|_2^2 \right) \\ &= \|z\|_2 - \frac{\|\tilde{x}\|_2^2 + 2v\lambda \|\tilde{x}\|_2^{2/3}}{2\|\tilde{x}\|_2} \\ &= \frac{1}{2} \|z\|_2 - \frac{2}{3} v\lambda \|\tilde{x}\|_2^{-1/3}, \end{aligned}$$

where the third equality holds since that \tilde{x} is proportional to z , and fourth equality follows from (3.17). Since both $\|z\|_2$ and $\|\tilde{x}\|_2$ are strictly increasing on $\|z\|_2$, $H(\|z\|_2)$ is also strictly increasing when $\|z\|_2 > 4(\frac{2}{9}v\lambda)^{3/4}$. Thus the unique solution of $H(\|z\|_2) = 0$ satisfies

$$\|z\|_2 \|\tilde{x}\|_2^{1/3} = \frac{4}{3} v\lambda,$$

and further, (3.17) implies that the solution of $H(\|z\|_2) = 0$ is

$$\|z\|_2 = 2 \left(\frac{2}{3} v\lambda \right)^{3/4}.$$

Therefore, we arrive at the formulae (3.9) and (3.10). ■

Remark 3.1. Note from (3.5), (3.6), (3.8), (3.9) and (3.11) that the solutions of the proximal optimization subproblems might not be unique when $Q_{p,q}(\tilde{x}) = Q_{p,q}(0)$. To avoid this obstacle in numerical computations, we select the solution $P_{p,q}(z) = 0$ whenever $Q_{p,q}(\tilde{x}) = Q_{p,q}(0)$, which achieves a more sparse solution, in the definition of the proximal operator to guarantee a unique update.

Remark 3.2. By Proposition 3.1, one sees that the proximal gradient method meets the group sparsity structure, since the components of each iterate within each group are likely to be either all zeros or all nonzeros. When $n_{\max} = 1$, the data do not form any group structure in the feature space, and the sparsity is achieved only on the individual feature level. In this case, the proximal operators $P_{2,1}(z)$, $P_{2,0}(z)$, and $P_{2,1/2}(z)$ and $P_{1,1/2}(z)$ reduce to the soft thresholding function in [18], the hard thresholding function in [6] and the half thresholding function in [57], respectively.

Remark 3.3. Proposition 3.1 presents the analytical solution of the proximal optimization subproblems (3.3) when $q = 0, 1/2, 2/3, 1$. However, in other cases, the analytical solution of (3.3) seems not available, since the algebraic equation (3.12) does not have an analytical solution (it is difficult to find an analytical solution for the algebraic equation whose order is larger than four). Thus, in the general cases of $q \in (0, 1)$, we alternatively use the Newton method to solve the nonlinear equation (3.12), which is the optimality condition of the proximal optimization subproblem. The numerical simulation in Figure 4 of Section 4 shows

that the Newton method works in solving the proximal optimization subproblems (3.3) for the general q , while the $\ell_{p,1/2}$ regularization is the best one among the $\ell_{p,q}$ regularizations for $q \in [0, 1]$.

3.2 Linear convergence rate

Recall that convergence results of PGM-GSO are from the references [3, 6, 7], saying that the generated sequence globally converges to a critical point or a global/local minimum of the $\ell_{p,q}$ regularization problem. However, the result on convergence rates of the proximal gradient method for solving lower-order regularization problems is still undiscovered. In this subsection, we will establish the linear convergence rate of PGM-GSO for the case $p = 1$ and $0 < q < 1$.

By virtue of the second-order necessary condition of (3.2), the following lemma provides a lower bound for nonzero groups of sequence $\{x^k\}$ generated by the PGM-GSO and shows that the index set of nonzero groups of $\{x^k\}$ maintains constant for large k .

Lemma 3.1. *Let $p = 1$, $0 < q < 1$ and $K = (v\lambda q(1-q))^{\frac{1}{2-q}}$. Let $\{x^k\}$ be a sequence generated by the PGM-GSO with $v < \frac{1}{2}\|A\|_2^{-2}$. Then the following statements hold:*

- (i) *for any i and k , if $x_{\mathcal{G}_i}^k \neq 0$, then $\|x_{\mathcal{G}_i}^k\|_1 \geq K$.*
- (ii) *x^k shares the same index set of nonzero groups for large k , that is, there exist $N \in \mathbb{N}$ and $\mathcal{I} \subseteq \{1, \dots, r\}$ such that*

$$\begin{cases} x_{\mathcal{G}_i}^k \neq 0, & i \in \mathcal{I}, \\ x_{\mathcal{G}_i}^k = 0, & i \notin \mathcal{I}, \end{cases} \quad \text{for all } k \geq N.$$

Proof. (i) For each group $x_{\mathcal{G}_i}^k$, by (3.2), one has that

$$x_{\mathcal{G}_i}^k \in \text{Arg} \min_{x \in \mathbb{R}^{n_i}} \left\{ \lambda \|x\|_1^q + \frac{1}{2v} \|x - z_{\mathcal{G}_i}^{k-1}\|_2^2 \right\}. \quad (3.23)$$

If $x_{\mathcal{G}_i}^k \neq 0$, we define $\mathcal{A}_i^k := \{j \in \mathcal{G}_i : x_j^k \neq 0\}$ and $a_i^k := |\mathcal{A}_i^k|$. Without loss of generality, we assume that the first a_i^k components of $x_{\mathcal{G}_i}^k$ are nonzeros. Then (3.23) implies that

$$x_{\mathcal{G}_i}^k \in \text{Arg} \min_{x \in \mathbb{R}^{a_i^k} \times \{0\}} \left\{ \lambda \|x\|_1^q + \frac{1}{2v} \|x - z_{\mathcal{G}_i}^{k-1}\|_2^2 \right\}. \quad (3.24)$$

The second-order necessary condition of (3.24) implies that

$$\frac{1}{v} I_i^k + \lambda q(q-1) M_i^k \succeq 0,$$

where I_i^k is the identity matrix in $\mathbb{R}^{a_i^k \times a_i^k}$ and $M_i^k = \|x_{\mathcal{A}_i^k}^k\|_1^{q-2} (\text{sign}(x_{\mathcal{A}_i^k}^k)) (\text{sign}(x_{\mathcal{A}_i^k}^k))^\top$. Let e be the first column of I_i^k . Therefore, we obtain that

$$\frac{1}{v} e^\top I_i^k e + \lambda q(q-1) e^\top M_i^k e \geq 0,$$

that is,

$$\frac{1}{v} + \lambda q(q-1) \|x_{\mathcal{A}_i^k}^k\|_1^{q-2} \geq 0.$$

Consequently, it implies that

$$\|x_{\mathcal{G}_i}^k\|_1 = \|x_{\mathcal{A}_i^k}^k\|_1 \geq (v\lambda q(1-q))^{\frac{1}{2-q}} = K.$$

Hence, it completes the proof of (i).

- (ii) Recall from Theorem 3.1 that $\{x^k\}$ converges to a critical point x^* . Then there exists $N \in \mathbb{N}$ such that $\|x^k - x^*\|_2 < \frac{K}{2\sqrt{n}}$, and thus,

$$\|x^{k+1} - x^k\|_2 \leq \|x^{k+1} - x^*\|_2 + \|x^k - x^*\|_2 < \frac{K}{\sqrt{n}}, \quad (3.25)$$

for any $k \geq N$. Proving by contradiction, without loss of generality, we assume that there exist $k \geq N$ and $i \in \{1, \dots, r\}$ such that $x_{\mathcal{G}_i}^{k+1} \neq 0$ and $x_{\mathcal{G}_i}^k = 0$. Then it follows from (i) that

$$\|x^{k+1} - x^k\|_2 \geq \frac{1}{\sqrt{n}} \|x^{k+1} - x^k\|_1 \geq \frac{1}{\sqrt{n}} \|x_{\mathcal{G}_i}^{k+1} - x_{\mathcal{G}_i}^k\|_1 \geq \frac{K}{\sqrt{n}}. \quad (3.26)$$

Hence we arrive at a contradiction from (3.25) and (3.26). The proof is complete. ■

The following lemma provides the first- and second-order conditions for a local minimum of $\ell_{1,q}$ regularization problem.

Lemma 3.2. *Let $p = 1$ and $0 < q < 1$. Assume that x^* is a local minimum of (1.4), and that any nonzero group of x^* is active; without loss of generality, we assume that x^* is of structure $x^* = (y^{*\top}, 0)^\top$ with*

$$y^* = (x_{\mathcal{G}_1}^{*\top}, \dots, x_{\mathcal{G}_S}^{*\top})^\top \text{ and } x_{\mathcal{G}_i}^* \neq \mathbf{0} \text{ for } i = 1, \dots, S. \quad (3.27)$$

Let $A = (B, D)$, where B is a submatrix corresponding to y^* , i.e., $B = (A_{\cdot j})$ with $j \in \{\mathcal{G}_i : i \in \mathcal{S}\}$ and $D = (A_{\cdot j})$ with $j \in \{\mathcal{G}_i : i \in \mathcal{S}^c\}$. Consider the following restricted problem

$$\min_{y \in \mathbb{R}^{n_{\mathbf{a}}}} f(y) + \varphi(y), \quad (3.28)$$

where $n_{\mathbf{a}} := \sum_{i \in \mathcal{S}} n_i$, and

$$f : \mathbb{R}^{n_{\mathbf{a}}} \rightarrow \mathbb{R} \quad \text{by} \quad f(y) := \|By - b\|_2^2 \quad \text{for any } y \in \mathbb{R}^{n_{\mathbf{a}}},$$

$$\varphi : \mathbb{R}^{n_{\mathbf{a}}} \rightarrow \mathbb{R} \quad \text{by} \quad \varphi(y) := \lambda \|y\|_{1,q}^q \quad \text{for any } y \in \mathbb{R}^{n_{\mathbf{a}}}.$$

Then the following statements are true:

(i) The following first- and second-order conditions hold

$$2B^\top(By^* - b) + \lambda q \begin{pmatrix} \|y_{\mathcal{G}_1}^*\|_1^{q-1} \text{sign}(y_{\mathcal{G}_1}^*) \\ \vdots \\ \|y_{\mathcal{G}_S}^*\|_1^{q-1} \text{sign}(y_{\mathcal{G}_S}^*) \end{pmatrix} = 0, \quad (3.29)$$

and

$$2B^\top B + \lambda q(q-1) \begin{pmatrix} M_1^* & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & M_S^* \end{pmatrix} \succ 0, \quad (3.30)$$

where $M_i^* = \|y_{\mathcal{G}_i}^*\|_1^{q-2} (\text{sign}(y_{\mathcal{G}_i}^*)) (\text{sign}(y_{\mathcal{G}_i}^*))^\top$.

(ii) The second-order growth condition holds at y^* for problem (3.28), that is, there exist $\varepsilon > 0$ and $\delta > 0$ such that

$$(f + \varphi)(y) \geq (f + \varphi)(y^*) + \varepsilon \|y - y^*\|_2^2 \quad \text{for any } y \in B(y^*, \delta). \quad (3.31)$$

Proof. (i) By (3.27), one has that $\varphi(\cdot)$ is smooth around y^* with its first- and second-derivatives being

$$\varphi'(y^*) = \lambda q \begin{pmatrix} \|y_{\mathcal{G}_1}^*\|_1^{q-1} \text{sign}(y_{\mathcal{G}_1}^*) \\ \vdots \\ \|y_{\mathcal{G}_S}^*\|_1^{q-1} \text{sign}(y_{\mathcal{G}_S}^*) \end{pmatrix},$$

and

$$\varphi''(y^*) = \lambda q(q-1) \begin{pmatrix} M_1^* & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & M_S^* \end{pmatrix};$$

hence $(f + \varphi)(\cdot)$ is also smooth around y^* . Therefore we obtain the following first- and second-order necessary conditions of (3.28)

$$f'(y^*) + \varphi'(y^*) = 0 \quad \text{and} \quad f''(y^*) + \varphi''(y^*) \succeq 0,$$

which are (3.29) and

$$2B^\top B + \lambda q(q-1) \begin{pmatrix} M_1^* & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & M_S^* \end{pmatrix} \succeq 0, \quad (3.32)$$

respectively. Proving by contradiction, we assume that (3.30) does not hold, i.e., there exists some $w \neq 0$ such that

$$2w^\top B^\top B w + \lambda q(q-1) \sum_{i=1}^S \left(\|y_{\mathcal{G}_i}^*\|_1^{q-2} \cdot \left(\sum_{j \in \mathcal{G}_i} w_j \text{sign}(y_j^*) \right)^2 \right) = 0.$$

Let $h : \mathbb{R} \rightarrow \mathbb{R}$ with $h(t) := \|B(y^* + tw) - b\|_2^2 + \lambda \|y^* + tw\|_p^p$. Clearly, $h(\cdot)$ has a local minimum at 0, and $h(\cdot)$ is smooth around 0 with its derivatives being

$$\begin{aligned} h'(0) &= 2w^\top B^\top (By^* - b) + \lambda q \sum_{i=1}^S \left(\|y_{\mathcal{G}_i}^*\|_1^{q-1} \cdot \sum_{j \in \mathcal{G}_i} w_j \text{sign}(y_j^*) \right) = 0, \\ h''(0) &= 2w^\top B^\top Bw + \lambda q(q-1) \sum_{i=1}^S \left(\|y_{\mathcal{G}_i}^*\|_1^{q-2} \cdot \left(\sum_{j \in \mathcal{G}_i} w_j \text{sign}(y_j^*) \right)^2 \right) = 0, \\ h^{(3)}(0) &= \lambda q(q-1)(q-2) \sum_{i=1}^S \left(\|y_{\mathcal{G}_i}^*\|_1^{q-3} \cdot \left(\sum_{j \in \mathcal{G}_i} w_j \text{sign}(y_j^*) \right)^3 \right) = 0, \end{aligned}$$

and

$$h^{(4)}(0) = \lambda q(q-1)(q-2)(q-3) \sum_{i=1}^S \left(\|y_{\mathcal{G}_i}^*\|_1^{q-4} \cdot \left(\sum_{j \in \mathcal{G}_i} w_j \text{sign}(y_j^*) \right)^4 \right) < 0. \quad (3.33)$$

However, it is clear that $h^{(4)}(0)$ must be nonnegative, which yields a contradiction to (3.33). Therefore, we proved (3.30).

- (ii) By the structure of y^* (cf. (3.27)), $\varphi(\cdot)$ is smooth around y^* , and thus, $(f + \varphi)(\cdot)$ is also smooth around y^* with its derivatives being

$$f'(y^*) + \varphi'(y^*) = 0 \quad \text{and} \quad f''(y^*) + \varphi''(y^*) \succ 0$$

(due to (3.29) and (3.30)). Hence the second-order growth condition (3.31) follows from [45, Theorem 13.24]. This completes the proof. ■

The key of convergence rate analysis of PGM-GSO is the descent of the functional $f + \varphi$ in each iteration step. The following lemma states some basic properties of active groups of sequence $\{x^k\}$ generated by the PGM-GSO.

Lemma 3.3. *Let $p = 1$ and $0 < q < 1$. Let $\{x^k\}$ be a sequence generated by the PGM-GSO with $v < \frac{1}{2}\|A\|_2^{-2}$, which converges to x^* (by Theorem 3.1). Let the assumptions and notations used in Lemma 3.2 be adopted. We further define*

$$\alpha := \|B\|_2^2, \quad L := 2\|A\|_2^2 \quad \text{and} \quad D_k := \varphi(y^k) - \varphi(y^{k+1}) + \langle f'(y^k), y^k - y^{k+1} \rangle.$$

Then there exist some $\delta > 0$ and $N \in \mathbb{N}$ such that the following inequalities hold for any $w \in B(y^, \delta)$ and any $k \geq N$:*

$$\varphi(w) - \varphi(y^{k+1}) + \langle f'(y^k), w - y^{k+1} \rangle \geq \frac{1}{v} \langle y^k - y^{k+1}, w - y^{k+1} \rangle - \alpha \|w - y^{k+1}\|_2^2, \quad (3.34)$$

$$D_k \geq \left(\frac{1}{v} - \alpha \right) \|y^k - y^{k+1}\|_2^2, \quad (3.35)$$

and

$$(f + \varphi)(y^{k+1}) \leq (f + \varphi)(y^k) - \left(1 - \frac{Lv}{2(1 - v\alpha)} \right) D_k. \quad (3.36)$$

Proof. By Lemma 3.1(ii) and the fact that $\{x^k\}$ converges to x^* , one has that x^k shares the same index set of nonzero groups with that of x^* for large k ; further by the structure of y^* (cf. (3.27)), we obtain that all components in nonzero groups of y^k are nonzero for large k . In another word, we have

$$\text{there exists } N \in \mathbb{N} \text{ such that } y_k \neq_{\mathbf{a}} 0 \text{ and } z^k = 0 \text{ for any } k \geq N; \quad (3.37)$$

hence $\varphi(\cdot)$ is smooth around y^k for any $k \geq N$

In view of PGM-GSO and the decomposition of $x = (y^\top, z^\top)^\top$, one has that

$$y^{k+1} \in \text{Arg min} \left\{ \lambda \varphi(y) + \frac{1}{2v} \left\| y - \left(y^k - v f'(y^k) \right) \right\|_2^2 \right\}. \quad (3.38)$$

The first-order necessary condition of (3.38) is

$$\varphi'(y^{k+1}) = \frac{1}{v} \left(y^k - v f'(y^k) - y^{k+1} \right). \quad (3.39)$$

Recall from (3.30) that $\varphi''(y^*) \succ -2B^\top B$. Since $\varphi(\cdot)$ is smooth around y^* , then there exists $\delta > 0$ such that $\varphi''(w) \succ -2B^\top B$ for any $w \in B(y^*, \delta)$. Noting that $\{y^k\}$ converges to y^* , without loss of generality, we assume that $\|y^k - y^*\| < \delta$ for any $k \geq N$ (otherwise, we can choose a larger N). Therefore, one has that $\varphi''(y^k) \succ -2B^\top B$ for any $k \geq N$. Then by Taylor expansion, we can assume without loss of generality that the following inequality holds for any $k \geq N$ and any $w \in B(y^*, \delta)$ (otherwise, we can choose a smaller δ):

$$\varphi(w) > \varphi(y^{k+1}) + \langle \varphi'(y^{k+1}), w - y^{k+1} \rangle - \alpha \|w - y^{k+1}\|_2^2.$$

Hence, by (3.39), it follows that

$$\varphi(w) - \varphi(y^{k+1}) > \frac{1}{v} \langle y^k - v f'(y^k) - y^{k+1}, w - y^{k+1} \rangle - \alpha \|w - y^{k+1}\|_2^2, \quad (3.40)$$

which is reduced to (3.34), and (3.35) follows by setting $w = y^k$ in (3.34). Furthermore, by the definition of $f(\cdot)$, it is of class $C_L^{1,1}$ and it follows from [4, Proposition A.24] that

$$\|f(y) - f(x) - f'(x)(y - x)\| \leq \frac{L}{2} \|y - x\|^2 \quad \text{for any } x, y.$$

Then, by the definition of D_k , it follows that

$$\begin{aligned} (f + \varphi)(y^{k+1}) - (f + \varphi)(y^k) + D_k &= f(y^{k+1}) - f(y^k) + \langle f'(y^k), y^k - y^{k+1} \rangle \\ &\leq \frac{L}{2} \|y^k - y^{k+1}\|_2^2 \\ &\leq \frac{Lv}{2(1 - v\alpha)} D_k, \end{aligned}$$

where the last inequality follows from (3.35), and thus, (3.36) is proved. \blacksquare

The main result of this subsection is presented as follows, where we prove the linear convergence rate of the PGM-GSO to a local minimum for the case $p = 1$ and $0 < q < 1$ under some mild assumptions.

Theorem 3.2. *Let $p = 1$ and $0 < q < 1$. Let $\{x^k\}$ be a sequence generated by the PGM-GSO with $v < \frac{1}{2}\|A\|_2^{-2}$. Then $\{x^k\}$ converges to a critical point x^* of (1.4). Further assume that x^* is a local minimum of (1.4), and that any nonzero group of x^* is active. Then there exist $N \in \mathbb{N}$, $C > 0$ and $\eta \in (0, 1)$ such that*

$$F(x^k) - F(x^*) \leq C\eta^k \quad \text{and} \quad \|x^k - x^*\|_2 \leq C\eta^k, \quad \text{for any } k \geq N. \quad (3.41)$$

Proof. The convergence of $\{x^k\}$ to a critical point x^* of (1.4) directly follows from Theorem 3.1. Let notations used in Lemma 3.2 be adopted, D_k , N and δ be defined as in Lemma 3.3, and let

$$r_k := F(x^k) - F(x^*).$$

Note in (3.37) that $y^k \neq_{\mathbf{a}} 0$ and $z^k = 0$ for any $k \geq N$. Thus

$$r_k = (f + \varphi)(y^k) - (f + \varphi)(y^*) \quad \text{for any } k \geq N.$$

It is trivial to see that $\varphi(\cdot)$ is smooth around y^* (as it is active), and that

$$\varphi''(y^*) = \lambda q(q-1) \begin{pmatrix} M_1^* & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & M_S^* \end{pmatrix} \prec 0, \quad f''(y^*) + \varphi''(y^*) \succ 0$$

(as shown in (3.30)). This shows that $\varphi(\cdot)$ is concave around y^* , while $(f + \varphi)(\cdot)$ is convex around y^* . Without loss of generality, we assume that $\varphi(\cdot)$ is concave and $(f + \varphi)(\cdot)$ is convex in $B(y^*, \delta)$ and that $y^k \in B(y^*, \delta)$ for any $k \geq N$ (since $\{y^k\}$ converges to y^*).

By the convexity of $(f + \varphi)(\cdot)$ in $B(y^*, \delta)$, it follows that for any $k \geq N$

$$\begin{aligned} r_k &= (f + \varphi)(y^k) - (f + \varphi)(y^*) \\ &\leq \langle f'(y^k) + \varphi'(y^k), y^k - y^* \rangle \\ &= \langle f'(y^k) + \varphi'(y^k), y^k - y^{k+1} \rangle + \langle f'(y^k) + \varphi'(y^k), y^{k+1} - y^* \rangle \\ &= D_k - \varphi(y^k) + \varphi(y^{k+1}) + \langle \varphi'(y^k), y^k - y^{k+1} \rangle + \langle f'(y^k) + \varphi'(y^k), y^{k+1} - y^* \rangle. \end{aligned} \quad (3.42)$$

Noting that $\varphi(\cdot)$ is concave in $B(y^*, \delta)$, it follows that

$$\varphi(y^k) - \varphi(y^{k+1}) \geq \langle \varphi'(y^k), y^k - y^{k+1} \rangle.$$

Consequently, (3.42) is reduced to

$$\begin{aligned} r_k &\leq D_k + \langle f'(y^k) + \varphi'(y^k), y^{k+1} - y^* \rangle \\ &= D_k + \langle \varphi'(y^k) - \varphi'(y^{k+1}), y^{k+1} - y^* \rangle + \langle f'(y^k) + \varphi'(y^{k+1}), y^{k+1} - y^* \rangle \\ &\leq D_k + \left(\frac{L_\varphi}{2} + \frac{1}{v} \right) \|y^k - y^{k+1}\|_2 \|y^{k+1} - y^*\|_2, \end{aligned} \quad (3.43)$$

where the last inequality follows from the smoothness of φ on $B(y^*, \delta)$ and (3.39), and L_φ is the Lipschitz constant of $\varphi'(\cdot)$ on $B(y^*, \delta)$. Let $\beta := 1 - \frac{Lv}{2(1-v\alpha)} \in (0, 1)$ (due to the assumption $v < \frac{1}{L}$). Then (3.36) is reduced to

$$r_k - r_{k+1} = (f + \varphi)(y^k) - (f + \varphi)(y^{k+1}) \geq \beta D_k > 0,$$

and thus, it follows from (3.43) and (3.35) that

$$\begin{aligned} \beta r_k &\leq \beta D_k + \beta \left(\frac{L_\varphi}{2} + \frac{1}{v} \right) \|y^k - y^{k+1}\|_2 \|y^{k+1} - y^*\|_2 \\ &\leq r_k - r_{k+1} + \beta \left(\frac{L_\varphi}{2} + \frac{1}{v} \right) \|y^{k+1} - y^*\|_2 \sqrt{\frac{v}{1-v\alpha}} D_k \\ &\leq r_k - r_{k+1} + \left(\frac{L_\varphi}{2} + \frac{1}{v} \right) \sqrt{\frac{v\beta}{1-v\alpha}} \|y^{k+1} - y^*\|_2 \sqrt{r_k - r_{k+1}}. \end{aligned} \quad (3.44)$$

Recall from Lemma 3.2(ii), there exists $c > 0$ such that

$$\|y - y^*\|_2^2 \leq c((f + \varphi)(y) - (f + \varphi)(y^*)) \quad \forall y \in B(y^*, \delta).$$

Thus, it follows that

$$\|y^{k+1} - y^*\|_2^2 \leq cr_{k+1} \leq cr_k \quad \text{for each } k \geq N. \quad (3.45)$$

Let $\epsilon := \frac{c}{\beta} \left(\frac{L_\varphi}{2} + \frac{1}{v} \right)^2$. By Young's inequality, (3.44) yields that

$$\begin{aligned} \beta r_k &\leq r_k - r_{k+1} + \frac{1}{2\epsilon} \|y^{k+1} - y^*\|_2^2 \left(\frac{L_\varphi}{2} + \frac{1}{v} \right)^2 + \frac{\epsilon v \beta}{2(1-v\alpha)} (r_k - r_{k+1}) \\ &\leq r_k - r_{k+1} + \frac{\beta}{2} r_k + \frac{cv}{2(1-v\alpha)} \left(\frac{L_\varphi}{2} + \frac{1}{v} \right)^2 (r_k - r_{k+1}). \end{aligned} \quad (3.46)$$

Let $\gamma := \frac{cv}{2(1-v\alpha)} \left(\frac{L_\varphi}{2} + \frac{1}{v} \right)^2 > 0$. Then (3.46) is reduced to

$$r_{k+1} \leq \frac{1 + \gamma - \frac{\beta}{2}}{1 + \gamma} r_k = \eta_1 r_k,$$

where $\eta_1 := \frac{1 + \gamma - \frac{\beta}{2}}{1 + \gamma} \in (0, 1)$. Thus, by letting $C_1 := r_N \eta_1^{-N}$, it follows that

$$r_k \leq \eta_1^{k-N} r_N = C_1 \eta_1^k \quad \text{for any } k \geq N.$$

By letting $\eta_2 = \sqrt{\eta_1}$ and $C_2 = \sqrt{C_1}$, it follows from (3.45) that

$$\|x^k - x^*\|_2 = \|y^k - y^*\|_2 \leq (cr_k)^{1/2} \leq C_2 \eta_2^k, \quad \text{for any } k \geq N.$$

By letting $C := \max\{C_1, C_2\}$ and $\eta := \max\{\eta_1, \eta_2\}$, we arrive at (3.48). The proof is complete. \blacksquare

Theorem 3.2 is an important theoretical result in that it establishes the linear convergence rate of proximal gradient method for solving the $\ell_{1,q}$ regularization problem under the assumption that any nonzero group of local minimum is an active group. Note that this assumption is satisfied automatically for the sparse optimization problem ($n_{max} = 1$). Hence, when $n_{max} = 1$, we obtain the linear convergence rate of proximal gradient method for solving ℓ_q regularization problem ($0 < q < 1$), which includes the iterative half thresholding algorithm ($q = 1/2$) proposed in [57] as a special case. This result is stated below, which we believe to the best of our knowledge that it is new.

Corollary 3.1. *Let $0 < q < 1$, and let $\{x^k\}$ be a sequence generated by the proximal gradient method for solving the following ℓ_q regularization problem*

$$\min_{x \in \mathbb{R}^n} F(x) := \|Ax - b\|_2^2 + \lambda \|x\|_q^q \quad (3.47)$$

with $v < \frac{1}{2}\|A\|_2^{-2}$. Then $\{x^k\}$ converges to a critical point x^ of (3.47). Further assume that x^* is a local minimum of (3.47). Then there exist $N \in \mathbb{N}$, $C > 0$ and $\eta \in (0, 1)$ such that*

$$F(x^k) - F(x^*) \leq C\eta^k \quad \text{and} \quad \|x^k - x^*\|_2 \leq C\eta^k, \quad \text{for any } k \geq N. \quad (3.48)$$

4 Numerical experiments

The purpose of this section is to carry out the numerical experiments of the proposed proximal gradient method for the $\ell_{p,q}$ regularization problem. We illustrate the performance of the PGM-GSO among different types of $\ell_{p,q}$ regularization, in particular, when $(p, q) = (2, 1), (2, 0), (2, 1/2), (1, 1/2), (2, 2/3)$ and $(1, 2/3)$, and compare them with several state-of-the-art algorithms, for both simulated data and real data in gene transcriptional regulation. All numerical experiments are implemented in MATLAB R2013b and executed on a personal desktop (Intel Core Duo E8500, 3.16 GHz, 4.00 GB of RAM).

4.1 Simulated data

In the numerical experiments on simulated data, the numerical data are generated as follows. We first randomly generate an i.i.d. Gaussian ensemble $A \in \mathbb{R}^{m \times n}$ satisfying $A^\top A = I$. Then we generate a group sparse solution $\bar{x} \in \mathbb{R}^n$ via randomly splitting its components into r groups and randomly picking k of them as active groups, whose entries are also randomly generated as i.i.d. Gaussian, while the remaining groups are all set as zeros. We generate the data b by the MATLAB script

$$b = A * \bar{x} + \text{sigma} * \text{randn}(m, 1),$$

where sigma is the standard deviation of additive Gaussian noise. The problem size is set to $n = 1024$ and $m = 256$, and we test on the noisy measurement data with $\text{sigma} = 0.1\%$. Assuming the group sparsity level S is predefined, the regularization parameter λ is iteratively

updated by obeying the rule: we set the iterative threshold to be the S -th largest value of $\|z_{\mathcal{G}_i}^k\|_2$ and solve the λ by virtue of Theorem 3.1.

For each given sparsity level, which is k/r , we randomly generate the data A , \bar{x} , b (as above) 500 times, run the algorithm, and average the 500 numerical results to illustrate the performance of the algorithm. We choose the stepsize $v = 1/2$ in all the testing. The two key criteria to characterize the performance are the relative error $\|x - \bar{x}\|_2 / \|\bar{x}\|_2$ and the successful recovery rate, where the recovery is defined as *success* when the relative error between the recovered data and the true data is smaller than 0.5%, otherwise, it is regarded as *failure*.

We carry out six experiments with the initial point $x_0 = 0$ (unless otherwise specified). In the first experiment, setting $r = 128$ (so group size $G = 1024/128 = 8$), we compare the convergence rate results and the successful recovery rates of the PGM-GSO with $(p, q) = (2, 1), (2, 0), (2, 1/2), (1, 1/2), (2, 2/3)$ and $(1, 2/3)$ for different sparsity levels. In Figure 2, (a), (b), and (c) illustrate the convergence rate results on sparsity level 1%, 5%, and 10%, respectively, while (d) plots the successful recovery rates on different sparsity levels. When the solution is of high sparse level, as shown in Figure 2(a), all $\ell_{p,q}$ regularization problems perform perfect and achieve a fast convergence rate. As demonstrated in Figure 2(b), when the sparsity level drops to 5%, $\ell_{p,1/2}$ and $\ell_{p,2/3}$ ($p = 1$ and 2) perform better and arrive at a more accurate level than $\ell_{2,1}$ and $\ell_{2,0}$. As illustrated in Figure 2(c), when the sparsity level is 10%, $\ell_{p,1/2}$ further outperforms $\ell_{p,2/3}$ ($p = 1$ or 2), and it surprises us that $\ell_{2,q}$ performs better and achieve a more accurate level than $\ell_{1,q}$ ($q = 1/2$ or $2/3$). From Figure 2(d), it is illustrated that $\ell_{p,1/2}$ achieves a better successful recovery rate than $\ell_{p,2/3}$ ($p = 1$ or 2), which outperforms $\ell_{2,0}$ and $\ell_{2,1}$. Moreover, we surprisingly see that $\ell_{2,q}$ also outperforms $\ell_{1,q}$ ($q = 1/2$ or $2/3$) on the successful recovery rate. In a word, $\ell_{2,1/2}$ performs as the best one of these six regularizations on both accuracy and robustness. In this experiment, we also note that the running times are at the same level, about 0.9 second per 500 iteration.

The second experiment is performed to show the sensitivity analysis on the group size ($G = 4, 8, 16, 32$) of the PGM-GSO with the six types of $\ell_{p,q}$ regularization. As shown in Figure 3, the six types of $\ell_{p,q}$ reach a higher successful recovery rate for the larger group size. We also note that the larger the group size, the shorter the running time.

The third experiment is implemented to study the variation of the PGM-GSO when varying the regularization order q (fix $p = 2$). Recall from Theorem 3.1, the analytical solution of the proximal optimization subproblems (3.3) can be obtained when $q = 0, 1/2, 2/3, 1$. However, in other cases, the analytical solution of (3.3) seems not available, and thus we apply the Newton method to solve the nonlinear equation (3.12), which is the optimality condition of the proximal optimization subproblem. Figure 4 shows the variation of successful recovery rates by decreasing the regularization order q from 1 to 0. It is illustrated that the PGM-GSO achieves the best successful recovery rate when $q = 1/2$, which arrives at the same conclusion as the first experiment. The farther the distance of q (in $[0, 1]$) from $1/2$, the lower the successful recovery rate.

The fourth experiment is to compare the PGM-GSO with several state-of-the-art algorithms in the field of sparse optimization, either convex or nonconvex algorithms, including

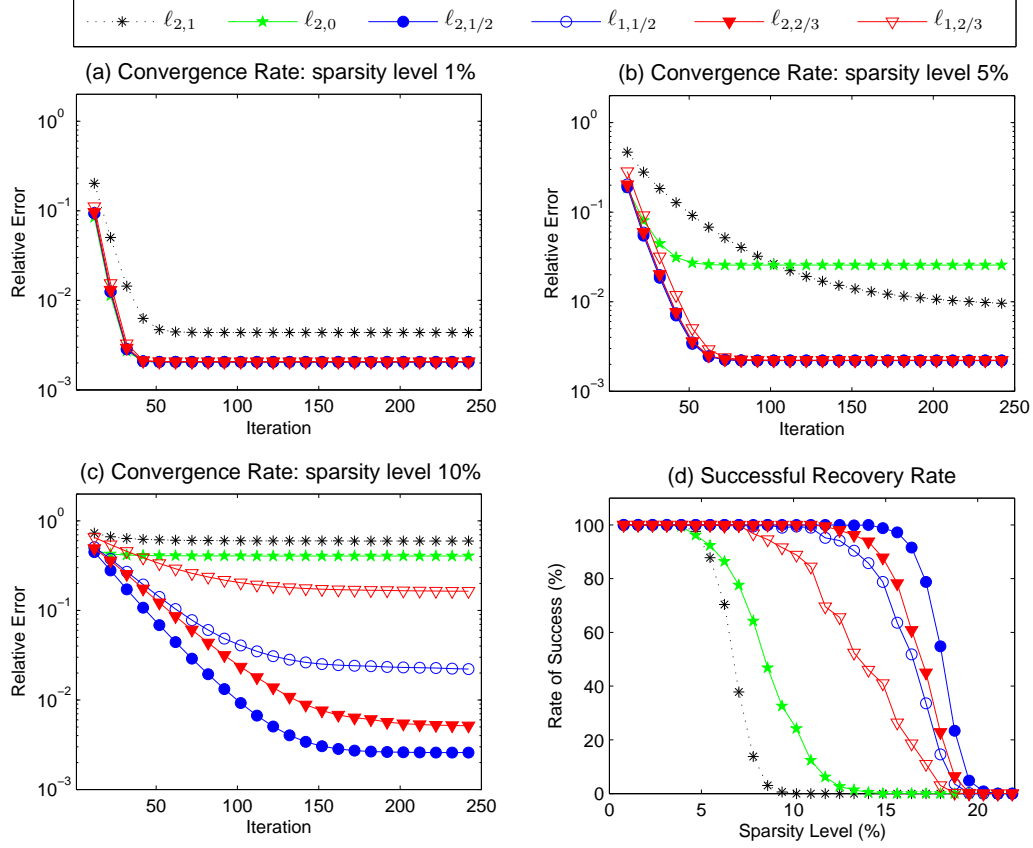


Figure 2: Convergence results and recovery rates for different sparsity levels.

ℓ_1 -Magic² [10], YALL1³ [19, 59], GBM⁴ [54], LqRecovery⁵ [26], HardTA⁶ [6, 8] and HalfTA⁷ [57]. Figure 5 demonstrates the successful recovery rates of these algorithms on different sparsity levels. It is indicated by Figure 5 that $\ell_{2,1/2}$ and $\ell_{2,2/3}$ can achieve the higher successful recovery rate than other algorithms, by exploiting the group sparsity structure.

Even though some global optimization method, such as the filled function method [28], can find the global solution of the lower-order regularization problem as in Example 2.2,

² ℓ_1 -Magic is a collection of MATLAB routines for solving the convex optimization programs central to compressive sampling, based on standard interior-point methods. The package is available at <http://users.ece.gatech.edu/~justin/l1magic/>.

³YALL1 (Your ALgorithm for L1) is a package of MATLAB solvers for the ℓ_1 sparse reconstruction, by virtue of the alternating direction method. The package is available at <http://yall1.blogs.rice.edu/>.

⁴GBM is a Gradient Based Method for solving the $\ell_{1/2}$ regularization problem. This algorithm is sensitive to the initial guess. Suggested by the authors, we choose the initial point as the solution obtained by the ℓ_1 -Magic.

⁵LqRecovery is an iterative algorithm for the ℓ_p norm minimization. The code is available at <http://www.math.drexel.edu/~foucart/software.htm>.

⁶HardTA is the iterative Hard Thresholding Algorithm, which is to solve the ℓ_0 regularization problem.

⁷HalfTA is the iterative Half Thresholding Algorithm, which is to solve the $\ell_{1/2}$ regularization problem.

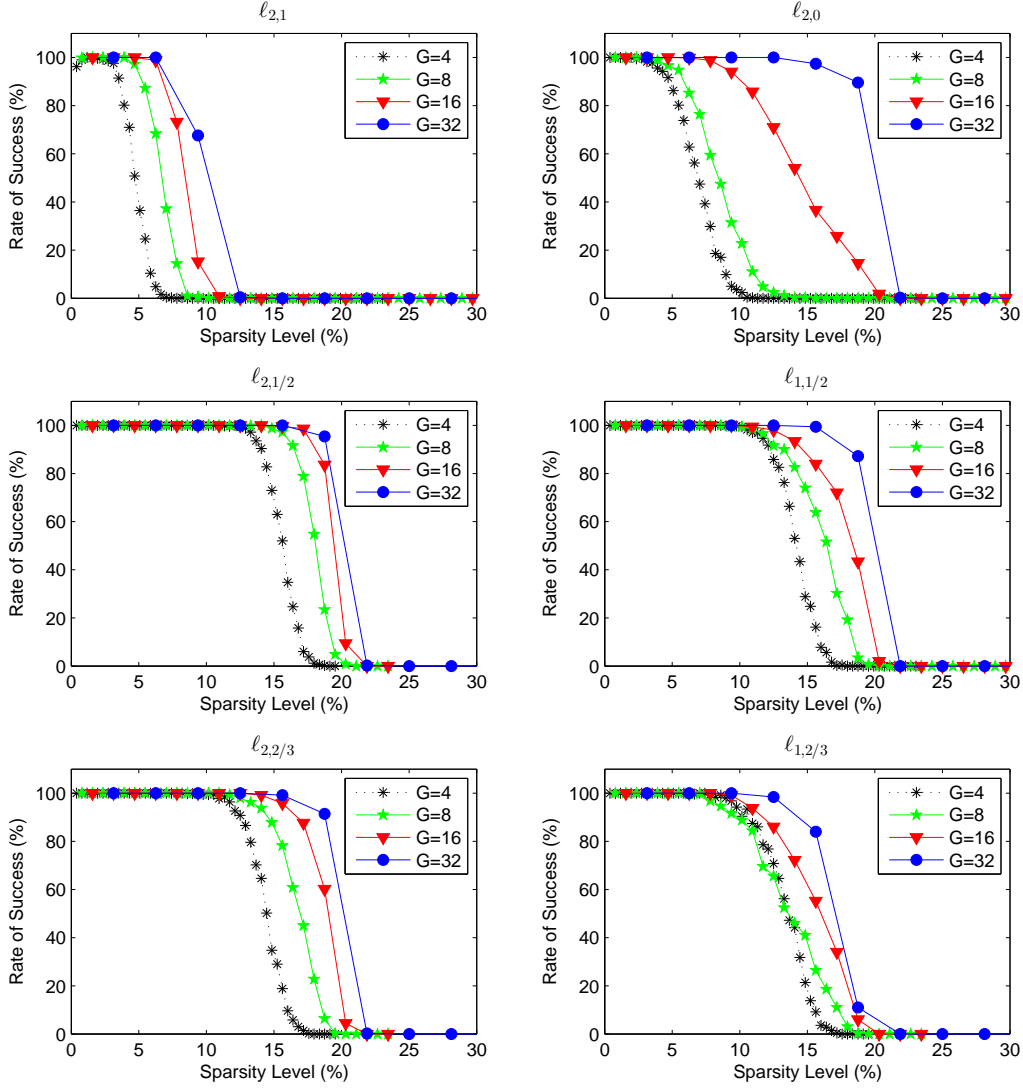


Figure 3: Sensitivity analysis on group size.

however, it does not work for the large-scale sparse optimization problems. Because, in the filled function method, all the directions need to be searched or compared in each iteration, which costs a large amount of time and hampers the efficiency for solving the large-scale problems.

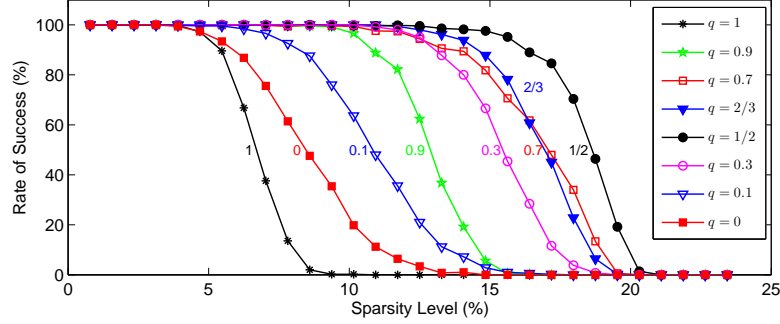


Figure 4: Variation of the PGM-GSO when varying the regularization order q .

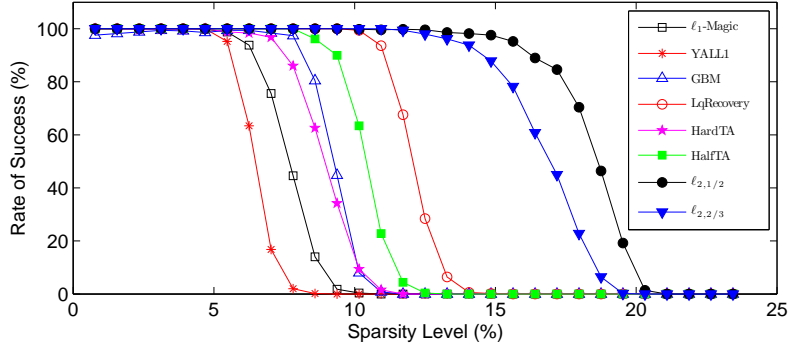


Figure 5: Comparison between the PGM-GSO and several state-of-the-art algorithms.

4.2 Real data in gene transcriptional regulation

Gene transcriptional regulation is the process that a combination of transcription factors (TFs) act in concert to control the transcription of the target genes. Inferring gene regulatory network from high-throughput genome-wide data is still a major challenge in systems biology, especially when the number of genes is large but the number of experimental samples is small. In large genomes, such as human and mouse, the complexity of gene regulatory system dramatically increases. Thousands of TFs combine in different ways to regulate tens of thousands target genes in various tissues or biological processes. However, only a few TFs collaborate and usually form complexes (i.e., groups of cooperative TFs) to control the expression of a specific gene in a specific cell type or developmental stage. Thus, the prevalence of TF complex makes the solution of gene regulatory network have a group structure, and the gene regulatory network inference in such large genomes becomes a group sparse optimization problem, which is to search a small number of TF complexes (or TFs) from a pool of thousands of TF complexes (or TFs) for each target gene based on the dependencies between the expression of TF complexes (or TFs) and the targets. Even though TFs often work in the form of complexes [56], and TF complexes are very important in the control of cell

identity and diseases [30], current methods to infer gene regulatory network usually consider each TF separately. To take the grouping information of TF complexes into consideration, we can apply the group sparse optimization to gene regulatory network inference with the prior knowledge of TF complexes as the pre-defined grouping.

4.2.1 Materials

Chromatin immunoprecipitation (ChIP) coupled with next generation sequencing (ChIP-seq) identifies *in vivo* active and cell-specific binding sites of a TF. They are commonly used to infer TF complexes recently. Thus, we manually collect ChIP-seq data in mouse embryonic stem cells (mESCs) (Table 1). Transcriptome is the gene expression profile of the whole genome that is measured by microarray or RNA-seq. The transcriptome data in mESCs for gene regulatory network inference are downloaded from Gene Expression Omnibus (GEO). 245 experiments under perturbations in mESC are collected from three papers [17, 40, 41]. Each experiment produced transcriptome data with or without overexpression or knockdown of a gene, in which the control and treatment have two replicates respectively. Gene expression fold changes between control samples and treatment samples of 12488 target genes in all experiments are log 2 transformed and form matrix $B \in \mathbb{R}^{245 \times 12488}$ (Figure ??A). The known TFs are collected from four TF databases, TRANSFAC, JASPAR, UniPROBE and TFCat, as well as literature. Let matrix $H \in \mathbb{R}^{245 \times 939}$ be made up of the expression profiles of 939 known TFs, and matrix $Z \in \mathbb{R}^{939 \times 12488}$ describe the connections between these TFs and targets. Then, the regulatory relationship between TFs and targets can be represented approximately by a linear system

$$HZ = B + \epsilon.$$

The TF-target connections defined by ChIP-seq data are converted into an initial matrix Z^0 (cf. [42]). Indeed, if TF i has a binding site around the gene j promoter within a defined distance (10 kbp), a non-zero number is assigned on Z_{ij}^0 as a prior value.

Now we add the grouping information (TF complexes) into this linear system. The TF complexes are inferred from ChIP-seq data (Table 1) via the method described in [29]. Let the group structure of Z be a matrix $W \in \mathbb{R}^{2257 \times 939}$ (actually, the number of groups is 1439), whose Moore-Penrose pseudoinverse [31] is denoted by W^+ . We further let $A := HW^+$ and $X := WZ$. Then the linear system can be converted into

$$AX = B + \epsilon,$$

where A denotes expression profiles of TF complexes, and X represents connections between TF complexes and targets (Figure ??A).

A literature-based golden standard (low-throughput golden standard) TF-target pair set from biological studies (Figure ??C), including 97 TF-target interactions between 23 TFs and 48 target genes, is downloaded from iScMiD (Integrated Stem Cell Molecular Interactions Database). Each TF-target pair in this golden standard dataset has been verified by biological experiments. Another more comprehensive golden standard mESC network is constructed

Table 1: ChIP-seq data for TF complex inference.

Factor	GEO accession	Pubmed ID	Factor	GEO accession	Pubmed ID
Atf7ip	GSE26680	-	Rad21	GSE24029	21589869
Atrx	GSE22162	21029860	Rbbp5	GSE22934	21477851
Cdx2	GSE16375	19796622	Rcor1	GSE27844	22297846
Chd4	GSE27844	22297846	Rest	GSE26680	-
Ctcf	GSE11431	18555785	Rest	GSE27844	22297846
Ctcf	GSE28247	21685913	Rnf2	GSE13084	18974828
Ctr9	GSE20530	20434984	Rnf2	GSE26680	-
Dpy30	GSE26136	21335234	Rnf2	GSE34518	22305566
E2f1	GSE11431	18555785	Setdb1	GSE17642	19884257
Ep300	GSE11431	18555785	Smad1	GSE11431	18555785
Ep300	GSE28247	21685913	Smad2	GSE23581	21731500
Esrrb	GSE11431	18555785	Smarca4	GSE14344	19279218
Ezh2	GSE13084	18974828	Smc1a	GSE22562	20720539
Ezh2	GSE18776	20064375	Smc3	GSE22562	20720539
Jarid2	GSE18776	20064375	Sox2	GSE11431	18555785
Jarid2	GSE19365	20075857	Stat3	GSE11431	18555785
Kdm1a	GSE27844	22297846	Supt5h	GSE20530	20434984
Kdm5a	GSE18776	20064375	Suz12	GSE11431	18555785
Klf4	GSE11431	18555785	Suz12	GSE13084	18974828
Lmnbl	GSE28247	21685913	Suz12	GSE18776	20064375
Med1	GSE22562	20720539	Suz12	GSE19365	20075857
Med12	GSE22562	20720539	Taf1	GSE30959	21884934
Myc	GSE11431	18555785	Taf3	GSE30959	21884934
Mycn	GSE11431	18555785	Tbp	GSE30959	21884934
Nanog	GSE11431	18555785	Tbx3	GSE19219	20139965
Nipbl	GSE22562	20720539	Tcfcp2l1	GSE11431	18555785
Nr5a2	GSE19019	20096661	Tet1	GSE26832	21451524
Pou5f1	GSE11431	18555785	Wdr5	GSE22934	21477851
Pou5f1	GSE22934	21477851	Whsc2	GSE20530	20434984
Prdm14	GSE25409	21183938	Zfx	GSE11431	18555785

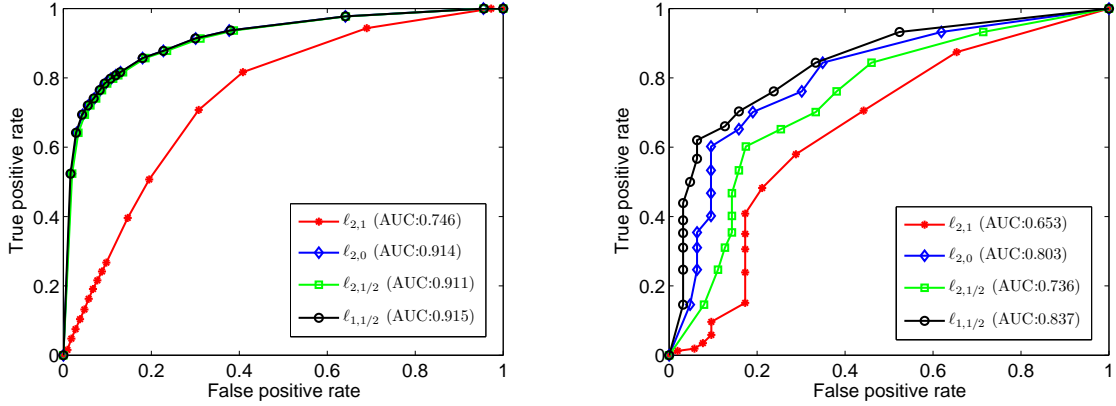
from high-throughput data (high-throughput golden standard) by ChIP-Array [43] using the methods described in [42]. It contains 40006 TF-target pairs between 13092 TFs or targets (Figure ??C). Basically, each TF-target pair in the network is evidenced by a cell-type specific binding site of the TF on the target’s promoter and the expression change of the target in the perturbation experiment of the TF, which is generally accepted as a true TF-target regulation. These two independent golden standards are both used to validate the accuracy of the inferred gene regulatory networks.

4.2.2 Numerical results

We apply and compare the PGM-GSO, starting from the initial matrix $X^0 := WZ^0$, to the gene regulatory network inference problem (Figure ??B). The area under the curve (AUC) of a receiver operating characteristic (ROC) curve is widely recognized as an important index of the overall classification performance of an algorithm; see [24]. Here, we apply AUC to evaluate the performance of the PGM-GSO with four types of $\ell_{p,q}$ regularization, $(p, q) = (2, 1), (2, 0), (2, 1/2)$ and $(1, 1/2)$. A series of numbers of predictive TF complexes (or TFs), denoted by k , from 1 to 100 (that is, the sparsity level varies from about 0.07% to 7%)

are tested. For each k and each pair of TF complex (or TF) i and target j , if the $X_{G_{ij}}^{(k)}$ is non-zero, this TF complex (or TF) is regarded as a potential regulator of this target in this test. In biological sense, we only concern about whether the true TF is predicted, but not the weight of this TF. We also expect that the TF complexes (or TFs) which are predicted in a higher sparsity level should be more important than those that are only reported in a lower sparsity level. Thus, when calculating the AUC, a score Score_{ij} is applied as the predictor for TF i on target j :

$$\text{Score}_{ij} := \begin{cases} \max_k \{1/k\}, & X_{ij}^{(k)} \neq 0, \\ 0, & \text{otherwise.} \end{cases}$$



(a) Evaluation with high-throughput golden standard.

(b) Evaluation with literature-based low-throughput golden standard.

Figure 6: ROC curves and AUCs of the PGM-GSO on mESC gene regulatory network inference.

Both high-throughput and low-throughput golden standards are used to draw the ROC curves of the PGM-GSO with four types of $\ell_{p,q}$ regularization in Figure 6 to compare their accuracy. When matched with the high-throughput golden standard, it is illustrated from Figure 6(a) that $\ell_{2,1/2}$, $\ell_{1,1/2}$ and $\ell_{2,0}$ perform almost the same (as indicated by the almost same AUC value), and significantly outperform $\ell_{2,1}$. With the low-throughput golden standard, it is demonstrated from Figure 6(b) that $\ell_{1,1/2}$ is slightly better than $\ell_{2,1/2}$ and $\ell_{2,0}$, and these three regularizations perform much better than $\ell_{2,1}$. These results are basically consistent with the results from simulated data. Since the golden standards we use here are obtained from real biological experiments, which are well-accepted as true TF-target regulations, the higher AUC, the more biologically accurate the result gene regulatory network is. Thus, our results indicate that the $\ell_{p,1/2}$ and $\ell_{p,0}$ regularizations are applicable to gene regulatory network inference in biological researches that study higher organisms but generate transcriptome data for only a small number of samples, which facilitates biologists to analyze gene regulation in a system level.

Acknowledgment. We are grateful to the four anonymous reviewers for their valuable suggestions and remarks which helped to improve the quality of the paper. We are also thankful to Professor Marc Teboulle for providing the reference [7] and the suggestion that the global convergence of the proximal gradient method for the nonconvex and nonsmooth composite optimization problem can be established by using the so-called Kurdyka-Łojasiewicz theory. Indeed, we show in Theorem 3.1 that the $\ell_{p,0}$ regularization problem globally converges to a local minimum and the $\ell_{p,q}$ ($0 < q < 1$) regularization problem globally converges to a critical point by virtue of the Kurdyka-Łojasiewicz theory.

References

- [1] H. Attouch, J. Bolte, P. Redont, and A. Soubeyran. Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the Kurdyka-Łojasiewicz inequality. *Mathematics of Operations Research*, 35(2):438–457, 2010.
- [2] F. R. Bach. Consistency of the group Lasso and multiple kernel learning. *Journal of Machine Learning Research*, 9:1179–1225, 2008.
- [3] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.
- [4] D. P. Bertsekas. *Nonlinear Programming*. Athena Scientific, Cambridge, 1999.
- [5] P. J. Bickel, Y. Ritov, and A. B. Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. *Annals of Statistics*, 37:1705–1732, 2009.
- [6] T. Blumensath and M. E. Davies. Iterative thresholding for sparse approximations. *Journal of Fourier Analysis and Applications*, 14:629–654, 2008.
- [7] J. Bolte, S. Sabach, and M. Teboulle. Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Mathematical Programming*, pages 459–494, 2013.
- [8] K. Bredies and D. Lorenz. Iterated hard shrinkage for minimization problems with sparsity constraints. *SIAM Journal on Scientific Computing*, 30(2):657–683, 2008.
- [9] F. Bunea, A. Tsybakov, and M. Wegkamp. Sparsity oracle inequalities for the Lasso. *Electronic Journal of Statistics*, 1:169–194, 2007.
- [10] E. Candès, J. Romberg, and T. Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Transactions on Information Theory*, 52(2):489–509, 2006.
- [11] E. Candès, J. K. Romberg, and T. Tao. Stable signal recovery from incomplete and inaccurate measurements. *Communications on Pure and Applied Mathematics*, 59(8):1207–1223, 2006.

- [12] E. Candès and T. Tao. Decoding by linear programming. *IEEE Transactions on Information Theory*, 51:4203–4215, 2005.
- [13] R. Chartrand and V. Staneva. Restricted isometry properties and nonconvex compressive sensing. *Inverse Problems*, 24:1–14, 2008.
- [14] S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM Review*, 43:129–159, 2001.
- [15] X. Chen, F. Xu, and Y. Ye. Lower bound theory of nonzero entries in solutions of ℓ_2 - ℓ_p minimization. *SIAM Journal on Scientific Computing*, 32(5):2832–2852, 2010.
- [16] P. Combettes and V. Wajs. Signal recovery by proximal forward-backward splitting. *Multiscale Modeling and Simulation*, 4(4):1168–1200, 2005.
- [17] L. S. Correa-Cerro, Y. Piao, A. A. Sharov, A. Nishiyama, J. S. Cadet, H. Yu, L. V. Sharova, L. Xin, H. G. Hoang, M. Thomas, Y. Qian, D. B. Dudekula, E. Meyers, B. Y. Binder, G. Mowrer, U. Bassey, D. L. Longo, D. Schlessinger, and M. S. Ko. Generation of mouse ES cell lines engineered for the forced induction of transcription factors. *Scientific Reports*, 1:167, 2011.
- [18] I. Daubechies, M. Defrise, and C. De Mol. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on Pure and Applied Mathematics*, 57:1413–1457, 2004.
- [19] W. Deng, W. Yin, and Y. Zhang. Group sparse optimization by alternating direction method. Technical report, Rice University, 2011.
- [20] D. L. Donoho. Compressed sensing. *IEEE Transactions on Information Theory*, 52(8):1289–1306, 2006.
- [21] D. L. Donoho. High-dimensional centrally symmetric polytopes with neighborliness proportional to dimension. *Discrete and Computational Geometry*, 35(4):617–652, 2006.
- [22] B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *Annals of Statistics*, 32:407–499, 2004.
- [23] J. Fan and R. Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96(456):1348–1360, 2001.
- [24] T. Fawcett. An introduction to ROC analysis. *Pattern Recognition Letters*, 27(8):861–874, 2006.
- [25] M. A. T. Figueiredo, R. D. Nowak, and S. J. Wright. Gradient projection for sparse reconstruction: Application to compressed sensing and other inverse problems. *IEEE Journal of Selected Topics in Signal Processing*, 1:586–597, 2007.

- [26] S. Foucart and M.-J. Lai. Sparsest solutions of underdetermined linear systems via ℓ_q -minimization for $0 < q \leq 1$. *Applied and Computational Harmonic Analysis*, 26(3):395–407, 2009.
- [27] D. Ge, X. Jiang, and Y. Ye. A note on complexity of L_p minimization. *Mathematical Programming, series B*, 129:285–299, 2011.
- [28] R. Ge. A filled function method for finding a global minimizer of a function of several variables. *Mathematical Programming*, 46(1-3):191–204, 1990.
- [29] E. G. Giannopoulou and O. Elemento. Inferring chromatin-bound protein complexes from genome-wide binding assays. *Genome Research*, 23(8):1295–1306, 2013.
- [30] D. Hnisz, B. J. Abraham, T. I. Lee, A. Lau, V. Saint-Andre, A. A. Sigova, H. A. Hoke, and R. A. Young. Super-enhancers in the control of cell identity and disease. *Cell*, 155(4):934–947, 2013.
- [31] R. A. Horn and C. R. Johnson. *Matrix Analysis*. Cambridge University Press, New York, 1985.
- [32] X. X. Huang and X. Q. Yang. A unified augmented Lagrangian approach to duality and exact penalization. *Mathematics of Operations Research*, 28(3):533–552, 2003.
- [33] Z. Luo, J. Pang, and D. Ralph. *Mathematical Programs with Equilibrium Constraints*. Cambridge University Press, Cambridge, 1996.
- [34] J. Mairal. Optimization with first-order surrogate functions. *International Conference on Machine Learning (ICML)*, 2013.
- [35] L. Meier, S. Van De Geer, and P. Bühlmann. The group Lasso for logistic regression. *Journal of the Royal Statistical Society, series B*, 70:53–71, 2008.
- [36] N. Meinshausen and B. Yu. Lasso-type recovery of sparse representations for high-dimensional data. *Annals of Statistics*, 37:246–270, 2009.
- [37] B. S. Mordukhovich. *Variational Analysis and Generalized Differentiation I: Basic Theory*. Springer, Berlin, 2006.
- [38] B. Natarajan. Sparse approximate solutions to linear systems. *SIAM Journal on Computing*, 24(2):227–234, 1995.
- [39] Y. Nesterov. Gradient methods for minimizing composite functions. *Mathematical Programming*, 140(1):125–161, 2013.
- [40] A. Nishiyama, A. A. Sharov, Y. Piao, M. Amano, T. Amano, H. G. Hoang, B. Y. Binder, R. Tapnio, U. Bassey, J. N. Malinou, L. S. Correa-Cerro, H. Yu, L. Xin, E. Meyers, M. Zalzman, Y. Nakatake, C. Stagg, L. Sharova, Y. Qian, D. Dudekula, S. Sheer, J. S.

- Cadet, T. Hirata, H. T. Yang, I. Goldberg, M. K. Evans, D. L. Longo, D. Schlessinger, and M. S. Ko. Systematic repression of transcription factors reveals limited patterns of gene expression changes in ES cells. *Scientific Reports*, 3:1390, 2013.
- [41] A. Nishiyama, L. Xin, A. A. Sharov, M. Thomas, G. Mowrer, E. Meyers, Y. Piao, S. Mehta, S. Yee, Y. Nakatake, C. Stagg, L. Sharova, L. S. Correa-Cerro, U. Bassey, H. Hoang, E. Kim, R. Tapnio, Y. Qian, D. Dudekula, M. Zalzman, M. Li, G. Falco, H. T. Yang, S. L. Lee, M. Monti, I. Stanghellini, M. N. Islam, R. Nagaraja, I. Goldberg, W. Wang, D. L. Longo, D. S. D, and M. S. Ko. Uncovering early response of gene regulatory networks in ESCs by systematic induction of transcription factors. *Cell Stem Cell*, 5:420–433, 2009.
- [42] J. Qin, Y. Hu, F. Xu, H. K. Yalamanchili, and J. Wang. Inferring gene regulatory networks by integrating ChIP-seq/chip and transcriptome data via LASSO-type regularization methods. *Methods*, 67(3):294–303, 2014.
- [43] J. Qin, M. J. Li, P. Wang, M. Q. Zhang, and J. Wang. ChIP-Array: Combinatory analysis of ChIP-seq/chip and microarray gene expression data to discover direct/indirect targets of a transcription factor. *Nucleic Acids Research*, pages W430–436, 2011.
- [44] M. Razaviyayn, M. Hong, and Z. Luo. A unified convergence analysis of block successive minimization methods for nonsmooth optimization. *SIAM Journal on Optimization*, 23(2):1126–1153, 2013.
- [45] R. T. Rockafellar and R. J.-B. Wets. *Variational Analysis*. Springer-Verlag, Berlin, 1998.
- [46] W. Rudin. *Principles of Mathematical Analysis*. McGraw-Hill, New York, 1976.
- [47] W. T. Short. Hyperbolic solution of the cubic equation. *National Mathematics Magazine*, 12(3):111–114, 1937.
- [48] R. Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society, series B*, 58:267–288, 1994.
- [49] P. Tseng and S. Yun. A coordinate gradient descent method for nonsmooth separable minimization. *Mathematical Programming*, 117(1-2):387–423, 2009.
- [50] M. Usman, C. Prieto, T. Schaeffter, and P. G. Batchelor. k - t Group sparse: A method for accelerating dynamic MRI. *Magnetic Resonance in Medicine*, 66(4):1163–1176, 2011.
- [51] S. A. van de Geer and P. Bühlmann. On the conditions used to prove oracle results for the Lasso. *Electronic Journal of Statistics*, 3:1360–1392, 2009.
- [52] E. van den Berg, M. Schmidt, M. Friedlander, and K. Murphy. Group sparsity via linear-time projection. Technical report, University of British Columbia, 2008.

- [53] S. J. Wright, R. D. Nowak, and M. A. T. Figueiredo. Sparse reconstruction by separable approximation. *IEEE Transactions on Signal Processing*, 57(7):2479–2493, 2009.
- [54] L. Wu, Z. Sun, and D.-H. Li. A gradient based method for the l_2 - $l_{1/2}$ minimization and application to compressive sensing. *Pacific Journal of Optimization*, 10(2):401–414, 2014.
- [55] L. Xiao and T. Zhang. A proximal-gradient homotopy method for the sparse least-squares problem. *SIAM Journal on Optimization*, 23(2):1062–1091, 2013.
- [56] D. Xie, A. P. Boyle, L. Wu, J. Zhai, T. Kawli, and M. Snyder. Dynamic trans-acting factor colocalization in human cells. *Cell*, 155(3):713–724, 2013.
- [57] Z. Xu, X. Chang, F. Xu, and H. Zhang. $L_{1/2}$ regularization: A thresholding representation theory and a fast solver. *IEEE Transactions on Neural Networks and Learning Systems*, 23:1013–1027, 2012.
- [58] H. Yang, Z. Xu, I. King, and M. R. Lyu. Online learning for group Lasso. In *International Conference on Machine Learning*, pages 1191–1198, 2010.
- [59] J. Yang and Y. Zhang. Alternating direction algorithms for ℓ_1 -problems in compressive sensing. *SIAM Journal on Scientific Computing*, 33(1):250–278, 2011.
- [60] X. Q. Yang and X. X. Huang. A nonlinear Lagrangian approach to constrained optimization problems. *SIAM Journal on Optimization*, 11(4):1119–1144, 2001.
- [61] M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of The Royal Statistical Society, series B*, 68:49–67, 2006.
- [62] T. Zhang. Some sharp performance bounds for least squares regression with L_1 regularization. *Annals of Statistics*, 37:2109–2144, 2009.
- [63] T. Zhang. Analysis of multi-stage convex relaxation for sparse regularization. *Journal of Machine Learning Research*, 11:1081–1107, 2010.